



UNITED STATES PATENT AND TRADEMARK OFFICE

UNITED STATES DEPARTMENT OF COMMERCE
United States Patent and Trademark Office
Address: COMMISSIONER FOR PATENTS
P.O. Box 1450
Alexandria, Virginia 22313-1450
www.uspto.gov

| APPLICATION NO. | FILING DATE | FIRST NAMED INVENTOR | ATTORNEY DOCKET NO. | CONFIRMATION NO. |
|-----------------|-------------|----------------------|---------------------|------------------|
| 09/936,174 | 12/06/2001 | David Naccache | 032326-166 | 9572 |

21839 7590 08/10/2005

BUCHANAN INGERSOLL PC
(INCLUDING BURNS, DOANE, SWECKER & MATHIS)
POST OFFICE BOX 1404
ALEXANDRIA, VA 22313-1404

EXAMINER

INGBERG, TODD D

| ART UNIT | PAPER NUMBER |
|----------|--------------|
|----------|--------------|

2193

DATE MAILED: 08/10/2005

Please find below and/or attached an Office communication concerning this application or proceeding.

Office Action Summary

Application No.

09/936,174

Applicant(s)

NACCACHE ET AL.

Examiner

Todd Ingberg

Art Unit

2124

- The MAILING DATE of this communication appears on the cover sheet with the correspondence address -
Period for Reply

A SHORTENED STATUTORY PERIOD FOR REPLY IS SET TO EXPIRE 3 MONTH(S) FROM THE MAILING DATE OF THIS COMMUNICATION.

- Extensions of time may be available under the provisions of 37 CFR 1.138(a). In no event, however, may a reply be timely filed after SIX (6) MONTHS from the mailing date of this communication.
- If the period for reply specified above is less than thirty (30) days, a reply within the statutory minimum of thirty (30) days will be considered timely.
- If NO period for reply is specified above, the maximum statutory period will apply and will expire SIX (6) MONTHS from the mailing date of this communication.
- Failure to reply within the set or extended period for reply will, by statute, cause the application to become ABANDONED (35 U.S.C. § 133). Any reply received by the Office later than three months after the mailing date of this communication, even if timely filed, may reduce any earned patent term adjustment. See 37 CFR 1.704(b).

Status

- 1) ☒ Responsive to communication(s) filed on 10 September 2001.
- 2a) ☐ This action is FINAL. 2b) ☒ This action is non-final.
- 3) ☐ Since this application is in condition for allowance except for formal matters, prosecution as to the merits is closed in accordance with the practice under *Ex parte Quayle*, 1935 C.D. 11, 453 O.G. 213.

Disposition of Claims

- 4) ☒ Claim(s) 1-26 and 28-31 is/are pending in the application.
- 4a) Of the above claim(s) 27 and 32 is/are withdrawn from consideration.
- 5) ☐ Claim(s) _____ is/are allowed.
- 6) ☒ Claim(s) 1-5, 7-13, 15, 18-20, 24-26 and 28-31 is/are rejected.
- 7) ☒ Claim(s) 6, 14, 16, 17, 21-23 is/are objected to.
- 8) ☐ Claim(s) _____ are subject to restriction and/or election requirement.

Application Papers

- 9) ☒ The specification is objected to by the Examiner.
- 10) ☒ The drawing(s) filed on 10 September 2001 is/are: a) ☐ accepted or b) ☒ objected to by the Examiner.
Applicant may not request that any objection to the drawing(s) be held in abeyance. See 37 CFR 1.85(a).
Replacement drawing sheet(s) including the correction is required if the drawing(s) is objected to. See 37 CFR 1.121(d).
- 11) ☒ The oath or declaration is objected to by the Examiner. Note the attached Office Action or form PTO-152.

Priority under 35 U.S.C. § 119

- 12) ☒ Acknowledgment is made of a claim for foreign priority under 35 U.S.C. § 119(a)-(d) or (f).
- a) ☒ All b) ☐ Some * c) ☐ None of:
1. ☒ Certified copies of the priority documents have been received.
 2. ☐ Certified copies of the priority documents have been received in Application No. _____.
 3. ☐ Copies of the certified copies of the priority documents have been received in this National Stage application from the International Bureau (PCT Rule 17.2(a)).

* See the attached detailed Office action for a list of the certified copies not received.

Attachment(s)

- 1) ☒ Notice of References Cited (PTO-892)
- 2) ☐ Notice of Draftsperson's Patent Drawing Review (PTO-948)
- 3) ☒ Information Disclosure Statement(s) (PTO-1449 or PTO/SB/08)
Paper No(s)/Mail Date 9/10/2001.
- 4) ☐ Interview Summary (PTO-413)
Paper No(s)/Mail Date. _____
- 5) ☐ Notice of Informal Patent Application (PTO-152)
- 6) ☐ Other: _____

Art Unit: 2124

DETAILED ACTION

Claims 1 – 26 and 28-31 have been examined.

In a Preliminary Amendment

Claims 1 – 26 and 28-31 were amended.

Claims 27 and 32 were cancelled.

Priority

1. Receipt is acknowledged of papers submitted under 35 U.S.C. 119(a)-(d), which papers have been placed of record in the file.

Information Disclosure Statement

2. The Information Disclosure Statement (IDS) filed December 6, 2001 has been considered. The reference in French could not be considered.

Oath/Declaration

3. Applicant has elected to use an outdated version of 37 CFR 1.56 "(as amended effective March 16, 1992)". Applicant should use the current form on the USPTO.GOV website when submitting a new Declaration.

Drawings

4. New corrected drawings in compliance with 37 CFR 1.121(d) are required in this application because they are fuzzy and hard to read and will not display properly in a U.S. Patent. Applicant is advised to employ the services of a competent patent draftsman outside the Office, as the U.S. Patent and Trademark Office no longer prepares new drawings. The corrected drawings are required in reply to the Office action to avoid abandonment of the application. The requirement for corrected drawings will not be held in abeyance.

Art Unit: 2124

5. Figure 1 should be designated by a legend such as --Prior Art-- because only that which is old is illustrated. See MPEP § 608.02(g). Corrected drawings in compliance with 37 CFR 1.121(d) are required in reply to the Office action to avoid abandonment of the application. The replacement sheet(s) should be labeled "Replacement Sheet" in the page header (as per 37 CFR 1.84(c)) so as not to obstruct any portion of the drawing figures. If the changes are not accepted by the examiner, the applicant will be notified and informed of any required corrective action in the next Office action. The objection to the drawings will not be held in abeyance.

Specification

6. The abstract of the disclosure is objected to because must be on a separate page.

Correction is required. See MPEP § 608.01(b).

7. Preliminary amendment of September 10, 2001 has been entered.

8. The disclosure is objected to because of the following informalities: The spelling of several words is not in the format for United States English, the European spelling of the following must be changed.

| <u>European Spelling</u> | <u>United States English Spelling</u> |
|--------------------------|---------------------------------------|
| "analysing" | analyzing |
| "analysed" | analyzed |
| "reinitialised" | reinitialized |
| "reinitialisation" | reinitialization |

Correction will benefit the searching of U.S. Patent literature.

Appropriate correction is required.

Art Unit: 2124

9. Page 5 of the Specification contains the acronym "FIBs", without the term being fully spelt out. On common meaning is "Secured hash standard, Federal Information Processing Standards Publication (FIPS) 180-1, May 1994". Clarification required with a change to the Specification.

10. The use of the trademark "JAVA" has been noted in this application. It should be capitalized wherever it appears and be accompanied by the generic terminology.

Although the use of trademarks is permissible in patent applications, the proprietary nature of the marks should be respected and every effort made to prevent their use in any manner which might adversely affect their validity as trademarks.

11. The title of the invention is not descriptive. A new title is required that is clearly indicative of the invention to which the claims are directed.

Claim Rejections - 35 USC § 112

12. The following is a quotation of the second paragraph of 35 U.S.C. 112:

The specification shall conclude with one or more claims particularly pointing out and distinctly claiming the subject matter which the applicant regards as his invention.

13. Claims 8 – 10 are rejected under 35 U.S.C. 112, second paragraph, as being indefinite for failing to particularly point out and distinctly claim the subject matter which applicant regards as the invention. The problem is the Applicant states the program to be monitored (DATA) . The focus of the claim language should the functionality of the monitor program and how it handles the varies condition presented by the input as it is processed. The Specification clearly supports what the Applicant is attempting to claim. This claim as written is indefinite. Dependent claims are also rejected merely because they are dependent on claim 8.

Art Unit: 2124

Claim 8

A method according to Claim 1 wherein, when the program to be monitored provides for at least one jump, the monitoring method is applied separately to sets of instructions in the program which do not include jumps between two successive instructions.

Claim Rejections - 35 USC § 102

14. The following is a quotation of the appropriate paragraphs of 35 U.S.C. 102 that form the basis for the rejections under this section made in this Office action:

A person shall be entitled to a patent unless –

(b) the invention was patented or described in a printed publication in this or a foreign country or in public use or on sale in this country, more than one year prior to the date of application for patent in the United States.

15. Claims 1 –5, 7-13, 15, 18-20, 24-26, 28, 30 and 31 are rejected under 35 U.S.C. 102(b) as being anticipated by USPN # 4,266,272 Berglund et al (IDS).

The environment of the invention JAVACARD is not claimed but is vastly different than the environment of the IDS reference

Claim Interpretation

The *control circuitry* in the reference IDS is performing the *monitor* function of the claimed invention.

Claim 1

IDS anticipates a method for monitoring progress with the execution of a linear sequence of instructions in a computer program (IDS, Abstract, control circuitry), comprising the steps of analysing the sequence of instructions transmitted to a processor intended to execute the program being monitored by extracting a data item from each instruction transmitted to the processor (IDS, Abstract, check word) and performing a calculation on said data item (IDS, Abstract, dynamically calculated), and verifying, the result of this analysis by comparing the result of said calculation to reference data (IDS, Abstract, local storage register vs. ALU), recorded with said program, wherein the reference data comprises a value pre-established so as to correspond to the result of the analysis produced during the monitoring method only if all the instructions in the sequence of instructions have actually been analysed during the running of the program (IDS, Abstract, control storage).

Claim Interpretation

The limitation “of a linear sequence of instructions” is not given patentable weight because it is dependent on the form of the input. Not part of the invention. It is treated as data.

Art Unit: 2124

Claim 2

A method according to Claim 1, wherein the verification of the result of the analysis is caused by an instruction placed at a predetermined location in the program to be monitored (as per claim 1 a register is a predetermined location), said instruction containing the reference data relating to a set of instructions whose correct execution is to be monitored (registers are inherently related to the instruction being processed).

Claim 3

A method according to Claim 1 wherein, when the instructions of the set of instructions to be monitored are in the form of a value, said analysis of the instructions is carried out by using these instructions as a numerical value. (Interpretation – all values are in binary format – this is inherent).

Claim 4

A method according to Claim 1, comprising the steps of:

- during the preparation of the program to be monitored (as per claim 1):
- incorporating, in at least one predetermined location in a sequence of instructions (as per claim 1) in the program, a reference value established according to a predetermined rule applied to identifiable data in each instruction to be monitored (as per claim 1, identification of words), and during the execution of the program to be monitored (as per claim 1):
- obtaining said identifiable data in each instruction received for execution (IDS, fetch col 9, 10-30),
- applying said predetermined rule to said identifiable data thus obtained in order to establish a verification value (as per claim 1), and
- verifying that this verification value actually corresponds to the reference value recorded with the program (as per claim 1).

Claim 5

A method according to Claim 1, further comprising a step of interrupting the flow of the program if the analysis reveals that the program being monitored has not been run as expected. (IDS, Figure #4, Result ERROR from result branch).

Claim 7

A method according to Claim 1 wherein the set of instructions to be monitored does not include jumps in its expected flow.

Claim Interpretation

The limitation “set of instructions to be monitored does not include jumps in its expected flow” is not given patentable weight because it is dependent on the form of the input. Not part of the invention. It is treated as data.

Claim 8

A method according to Claim 1 wherein, when the program to be monitored provides for at least one jump, the monitoring method is applied separately to sets of instructions in the program which do not include jumps between two successive instructions (IDS, col 9, lines 10 – 40).

Art Unit: 2124

Claim Interpretation

The limitation "program to be monitored provides for at least one jump" is not given patentable weight because it is dependent on the form of the input. Not part of the invention. It is treated as data.

Claim 9

A method according to Claim 8, wherein, when the program to be monitored includes an instruction for a jump dependent on the manipulated data, the monitoring method is implemented separately for a set of instructions which precedes the jump, and for at least one set of instructions which follows said jump. As per claim 8.

Claim 10

A method according to Claim 9, wherein, for a set of instructions providing for a jump, an instruction which controls this jump is integrated in said set of instructions for the purpose of obtaining a verification value for this set of instructions before executing the jump instruction. as per claim 8.

Claim 11

A method according to Claim 1 wherein the analysis is reinitialised before each new monitoring of a sequence of instructions to be monitored. (IDS, cycle and incrementer, col 3, lines 40 – 60)

Claim 12

A method according to Claim 11, wherein the reinitialisation of the analysis of each new monitoring includes the step of erasing or replacing a verification value obtained during a previous analysis. As per claim 11 depending on cycle determination.

Claim 13

A method according to Claim 11 wherein the reinitialisation of the monitoring analysis is controlled by the software itself. (Interpretation – the control circuitry and software being executed has a functional relationship – This is deemed inherent and related to Examiner's note above)

Claim 15

A method according to Claim 1 wherein the analysis includes the step of calculating, for each instruction under consideration following a previous instruction, the result of an operation on both a value obtained of the instruction in question and the result obtained by the same operation performed on the previous instruction. As per claim 1.

Claim 18

A method according to Claim 1 wherein the analysis includes the step of obtaining a comparison value by calculating successive intermediate values as the data of the respective instructions are obtained. (IDS, Abstract, last sentence words is plural).

Claim 19

Art Unit: 2124

A method according to Claim 1 wherein the analysis comprises a step of saving each data item necessary for verification, obtained from instructions in the set of instructions to be monitored as they are executed, and performing a calculation of a verification value from these data only at the necessary time, once all the necessary data have been obtained. (as per claim 1 and details of fetch col 3, lines 10-30).

Claim 20

IDS anticipates a device for monitoring progress with the execution of a series of instructions of a computer program, comprising means for analysing the sequence of instructions transmitted to the processor intended to execute the program being monitored by extracting a data item from each instruction transmitted to the processor and performing a calculation on said data item, and means for verifying the result of this analysis by comparing the result of said calculation to reference data recorded with said program, wherein the reference data comprises a value pre-established so as to correspond to the result of the analysis produced during monitoring only if all the instructions in the sequence of instructions have actually been analysed during the running of the program. As per claim 1.

Claim Interpretation

The limitation "a series of instructions of a computer program" is not given patentable weight because it is dependent on the form of the input. Not part of the invention. It is treated as data.

Claim 24

A device according to Claim 20 that is integrated into a programmed device containing said program to be monitored. (IDS, Abstract, Control Circuitry).

Claim 25

A device according to Claim 20 that is integrated into a program execution device. (IDS, Abstract, Control Circuitry).

Claim 26

IDS anticipates a program execution device that executes a series of instructions of a computer program, comprising means for analysing the sequence of instructions transmitted for execution by extracting a data item from each instruction and performing a calculation on said data item, and means for verifying the result of this analysis by comparing the result of said calculation to reference data recorded with the program to be monitored, wherein the reference data comprises a value pre-established so as to correspond to the result of the analysis produced during monitoring only if all the instructions in the sequence of instructions have actually been analysed during the running of the program. As per claim 1.

Claim Interpretation

A. The limitation "a series of instructions of a computer program" is not given patentable weight because it is dependent on the form of the input. Not part of the invention. It is treated as data.

B. In a similar fashion. The limitation "correspond to the result of the analysis produced during monitoring only if all the instructions in the sequence of instructions have actually been analysed during the running of the program" can be dependent on the input. If the program is only a few statements which all statements are to execute the claim limitations are input dependent. The

Art Unit: 2124

claim limitations outside the fact a monitor function is present have not performed a non required step to distinguish it from a monitor.

Claim 28

IDS anticipates a programmed device containing a series of recorded instructions and a fixed memory containing reference data pre-established as a function of data contained in said instructions for analysis and verification of the sequence of instructions, wherein the reference data comprises a value pre-established so as to correspond to the result of the analysis produced during monitoring only if all the instructions in the sequence of instructions have actually been analyzed during the running of the program. as per claim 1.

Claim Interpretation

A. The limitation "a series of recorded instructions" is not given patentable weight because it is dependent on the form of the input. Not part of the invention. It is treated as data.

Claim 30

A device according to Claim 28 wherein the reference data are recorded in the form of a prewired value or values fixed in memory. (IDS , Abstract, last sentence).

Claim Interpretation

The presence of the OR in the limitations. the Examiner elects to reject the underlined limitation above.

Claim 31

A device for programming a programmed device according to Claim 28, comprising means for entering, in at least one predetermined location in a sequence of instructions in the program, a reference value calculated according to a preestablished mode from data included in each instruction in a set of instructions whose execution is to be monitored. As per claim 1.

Claim Rejections - 35 USC § 103

16. The following is a quotation of 35 U.S.C. 103(a) which forms the basis for all

obviousness rejections set forth in this Office action:

(a) A patent may not be obtained though the invention is not identically disclosed or described as set forth in section 102 of this title, if the differences between the subject matter sought to be patented and the prior art are such that the subject matter as a whole would have been obvious at the time the invention was made to a person having ordinary skill in the art to which said subject matter pertains. Patentability shall not be negated by the manner in which the invention was made.

17. Claim 29 is rejected under 35 U.S.C. 103(a) as being unpatentable over IDS in view of

USPN # 6,402,028.

Claim 29

Art Unit: 2124

IDS teaches a device but IDS does not teach the device is a smart card. according to Claim 28, wherein said device is a smart card. USPN 6,402,028 teaches the production of Smart Cards where the logic is on the card. therefore, it would have been obvious to one of ordinary skill in the art at the time of invention, to combine IDS with 6,402,028 because logic control for Smart cards makes Smart Cards more reliable.

Allowable Subject Matter

18. Claims 6, 14, 16, 17 and 21 – 23 are objected to as being dependent upon a rejected base claim, but would be allowable if rewritten in independent form including all of the limitations of the base claim and any intervening claims. The bold and underlined limitations below indicate limitations not found in the prior art of record.

Claim 6

A method according to Claim 1, further comprising an invalidation step for future use of the device comprising the monitored program if **said analysis reveals a predetermined number of times that the program being monitored has not run in the expected manner.**

Claim 14

A method according to Claim 1 wherein the analysis produces a verification value obtained as the last value **in a series of values which is made to change successively with the analysis of each of the analysed instructions of the set of instructions, thus making it possible to contain an internal state of the running of the monitoring method and to follow its changes.**

Claim 16

A method according to Claim 1 wherein the analysis includes the step **of recursively applying a hash function to values obtained of each monitored instruction,** starting from a last initialisation performed.

Claim 17

A method according to Claim 1 wherein the analysis includes the step of making a verification value change by **performing a redundancy calculation on all the operating codes and the addresses executed since the last initialisation was carried out.**

Claim 21

A device according to Claim 20, further including a register for recording intermediate results in **a calculation in a chain carried out by the analysis means in order to obtain a verification value.**

Art Unit: 2124

Claim 22

A device according to Claim 21, further comprising means for recording a predetermined value or resetting the register under the control of an instruction transmitted during the execution of a program to be monitored. (Dependent on claim 21)

Claim 23

A device according to Claim 20, further comprising means for counting the number of unexpected events in the program being monitored, as determined by the analysis means, and means for invalidating the future use of the program to be monitored if this number reaches a predetermined threshold.

Conclusion

19. The prior art made of record and not relied upon is considered pertinent to applicant's disclosure.

US Patent Literature

- A. 6,402,028 – Deals with mass production of Smart Cards Column 4 covers JAVACARD technology.
- B. 6,668,325 – Employs an obfuscation technique on a section of code. Environment is distributed.
- C. 5,974,549 – Monitor is implemented via Dynamic Link Library (DLL).
- D. 6,546,546 – Appears to be dependent on the extensible operating system disclosed (PARAMETIUM).
- E. 6,092,120 – Focus on class loaders.
- F. 6,327,700 – Based on Profile data.
- G. 6,557,168 – Monitor is included at class level not at low level as per disclosed invention.

Art Unit: 2124

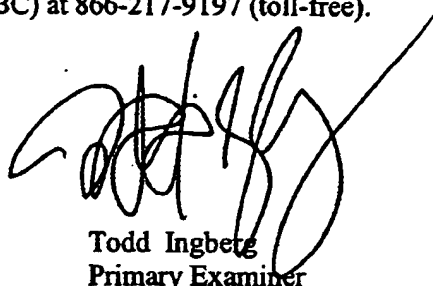
H. 6,275,938 – The monitor environment runs at operating system level not processor level as disclosed invention.

Correspondence

20. Any inquiry concerning this communication or earlier communications from the examiner should be directed to Todd Ingberg whose telephone number is (571) 272-3723. The examiner can normally be reached on during the work week..

If attempts to reach the examiner by telephone are unsuccessful, the examiner's supervisor, Kakali Chaki can be reached on (571) 272-3719. The fax phone number for the organization where this application or proceeding is assigned is 703-872-9306.

Information regarding the status of an application may be obtained from the Patent Application Information Retrieval (PAIR) system. Status information for published applications may be obtained from either Private PAIR or Public PAIR. Status information for unpublished applications is available through Private PAIR only. For more information about the PAIR system, see <http://pair-direct.uspto.gov>. Should you have questions on access to the Private PAIR system, contact the Electronic Business Center (EBC) at 866-217-9197 (toll-free).



Todd Ingberg
Primary Examiner
Art Unit 2124

TI

INFORMATION DISCLOSURE STATEMENT BY APPLICANT

GROUP
~~Unassigned~~ 2124

[illegible]

(05/01)

| | | | |
|-----------------------------------|---------------------------------------|---|-------------|
| Notice of References Cited | Application/Control No. 09/936,174 | Applicant(s)/Patent Under Reexamination NACCACHE ET AL. | |
| | Examiner Todd Ingberg | Art Unit 2124 | Page 1 of 3 |

U.S. PATENT DOCUMENTS

| * | | Document Number Country Code-Number-Kind Code | Date MM-YYYY | Name | Classification |
|---|---|--|-----------------|---------------------------|----------------|
| | A | US-6,862,684 | 03-2005 | DiGiorgio, Rinaldo | 713/163 |
| | B | US-6,402,028 | 06-2002 | Graham et al. | 235/380 |
| | C | US-6,418,420 | 07-2002 | DiGiorgio et al. | 705/40 |
| | D | US-6,615,264 | 09-2003 | Stoltz et al. | 709/227 |
| | E | US-6,581,206 | 06-2003 | Chen, Zhiqun | 717/143 |
| | F | US-6,023,764 | 02-2000 | Curtis, Bryce Allen | 713/200 |
| | G | US-5,983,348 | 11-1999 | Ji, Shuang | 713/200 |
| | H | US-6,092,120 | 07-2000 | Swaminathan et al. | 709/247 |
| | I | US-5,974,549 | 10-1999 | Golan, Gilad | 713/200 |
| | J | US-6,802,054 | 10-2004 | Faraj, Mazen | 717/128 |
| | K | US-6,546,546 | 04-2003 | Van Doorn, Leendert Peter | 717/114 |
| | L | US-6,557,168 | 04-2003 | Czajkowski, Grzegorz J. | 717/151 |
| | M | US-6,199,181 | 03-2001 | Rechef et al. | 714/38 |

FOREIGN PATENT DOCUMENTS

| * | | Document Number Country Code-Number-Kind Code | Date MM-YYYY | Country | Name | Classification |
|---|---|--|-----------------|---------|------|----------------|
| | N | | | | | |
| | O | | | | | |
| | P | | | | | |
| | Q | | | | | |
| | R | | | | | |
| | S | | | | | |
| | T | | | | | |

NON-PATENT DOCUMENTS

| * | | Include as applicable: Author, Title Date, Publisher, Edition or Volume, Pertinent Pages) |
|---|---|---|
| | U | JAVA 2 Complete, SYBEX, Steven Holzner, pages 3-28, 1999 |
| | V | "A New Public Key Cryptosystem Based on Higher Residues", David Naccache et al, ACM 1998, pages 59 - 66 |
| | W | "Twin Signatures:An Alternative to the Hash-and-Sign Paradigm", David Naccache et al, ACM 2001, pages 20 - 27 |
| | X | "Batch Exponentiation A Fast DLP-based Signature Generation Strategy", David M'Raihi and David Naccache, ACM, 1996, pages 58 - 61 |

*A copy of this reference is not being furnished with this Office action. (See MPEP § 707.05(a).)
Dates in MM-YYYY format are publication dates. Classifications may be US or foreign.

Notice of References Cited

Application/Control No.

09/936,174

Applicant(s)/Patent Under
Reexamination
NACCACHE ET AL.

Examiner

Todd Ingberg

Art Unit

2124

Page 2 of 3

U.S. PATENT DOCUMENTS

| * | | Document Number Country Code-Number-Kind Code | Date MM-YYYY | Name | Classification |
|---|---|--|-----------------|-----------------|----------------|
| | A | US-6,275,938 | 08-2001 | Bond et al. | 713/200 |
| | B | US-6,327,700 | 12-2001 | Chen et al. | 717/127 |
| | C | US-6,668,325 | 12-2003 | Collberg et al. | 713/194 |
| | D | US-6,510,352 | 01-2003 | Badavas et al. | 700/19 |
| | E | US-5,991,414 | 11-1999 | Garay et al. | 713/165 |
| | F | US-5,933,498 | 08-1999 | Schneck et al. | 705/54 |
| | G | US-6,314,409 | 11-2001 | Schneck et al. | 705/54 |
| | H | US-6,859,533 | 02-2005 | Wang et al. | 380/28 |
| | I | US-5,347,581 | 09-1994 | Naccache et al. | 380/30 |
| | J | US-5,452,357 | 09-1995 | Naccache, David | 713/172 |
| | K | US-6,698,662 | 03-2004 | Feyt et al. | 235/492 |
| | L | US-2003/0079127 | 04-2003 | Bidan et al. | 713/172 |
| | M | US-2004/0088555 | 05-2004 | Girard et al. | 713/192 |

FOREIGN PATENT DOCUMENTS

| * | | Document Number Country Code-Number-Kind Code | Date MM-YYYY | Country | Name | Classification |
|---|---|--|-----------------|---------|------|----------------|
| | N | | | | | |
| | O | | | | | |
| | P | | | | | |
| | Q | | | | | |
| | R | | | | | |
| | S | | | | | |
| | T | | | | | |

NON-PATENT DOCUMENTS

| * | | Include as applicable: Author, Title Date, Publisher, Edition or Volume, Pertinent Pages) |
|---|---|--|
| | U | "PicoDBMS:Scaling Down Database Technique For The Smartcard", Philippe Pucheral et al, VLDL Journal, 2001, pages 120-132 |
| | V | "Implementation for Coalesced Hashing", Jeffrey Scott Vitter Brown University, ACM, December 1982, pages 911-926 |
| | W | "Optimal Arrangement of Keys in Hash Table", Ronald L. Rivest, ACM, April 1978, pages 200-209 |
| | X | "Code Optimization Techniques for Embedded DSP Microprocessors", Stan Liao et al, ACM, 1995, 6 pages |

*A copy of this reference is not being furnished with this Office action. (See MPEP § 707.05(a).)
Dates in MM-YYYY format are publication dates. Classifications may be US or foreign.

| | | | |
|-----------------------------------|---------------------------------------|---|-------------|
| Notice of References Cited | Application/Control No. 09/936,174 | Applicant(s)/Patent Under Reexamination NACCACHE ET AL. | |
| | Examiner Todd Ingberg | Art Unit 2124 | Page 3 of 3 |

U.S. PATENT DOCUMENTS

| * | | Document Number Country Code-Number-Kind Code | Date MM-YYYY | Name | Classification |
|---|---|--|-----------------|-----------------------|----------------|
| | A | US-2002/0174309 | 11-2002 | Naccache et al. | 711/163 |
| | B | US-2003/0188170 | 10-2003 | Bidan et al. | 713/182 |
| | C | US-6,279,123 | 08-2001 | Mulrooney, Timothy J. | 714/35 |
| | D | US-6,507,904 | 01-2003 | Ellison et al. | 712/229 |
| | E | US-6,065,108 | 05-2000 | Tremblay et al. | 712/201 |
| | F | US-6,021,469 | 02-2000 | Tremblay et al. | 711/125 |
| | G | US-6,014,723 | 01-2000 | Tremblay et al. | 711/1 |
| | H | US- | | | |
| | I | US- | | | |
| | J | US- | | | |
| | K | US- | | | |
| | L | US- | | | |
| | M | US- | | | |

FOREIGN PATENT DOCUMENTS

| * | | Document Number Country Code-Number-Kind Code | Date MM-YYYY | Country | Name | Classification |
|---|---|--|-----------------|---------|------|----------------|
| | N | | | | | |
| | O | | | | | |
| | P | | | | | |
| | Q | | | | | |
| | R | | | | | |
| | S | | | | | |
| | T | | | | | |

NON-PATENT DOCUMENTS

| * | | Include as applicable: Author, Title Date, Publisher, Edition or Volume, Pertinent Pages) |
|---|---|--|
| | U | "Fundamental Technique for Order Optimization", David Simmen et al, ACM, 1996, pages 57 - 67 |
| | V | |
| | W | |
| | X | |

*A copy of this reference is not being furnished with this Office action. (See MPEP § 707.05(a).)
Dates in MM-YYYY format are publication dates. Classifications may be US or foreign.

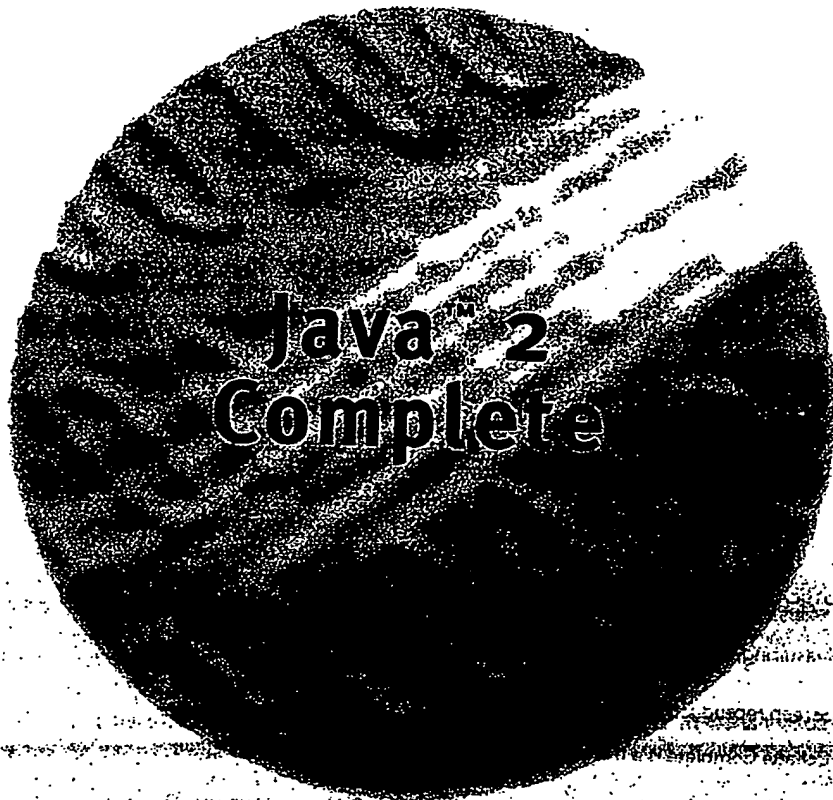
World's #1 Java Book Value

1,000 Pages
ONLY \$19.99! U.S.

JavaTM 2 COMPLETE

- Learn Java Essentials
- Build Interactive Web Applications with the JDK and Work with Java Objects
- Create Sophisticated GUIs with JFC Swing Components and the 2D API
- Learn About Java Beans





SYBEX SAN FRANCISCO ► PARIS ► DÜSSELDORF ► SOEST ► LONDON

Associate Publisher: Gary Masters

Contracts and Licensing Manager: Kristine O'Callaghan

Acquisitions & Developmental Editors: Denise Santoro and Maureen Adams

Project Editor: Gemma O'Sullivan

Compilation Editor: Suzanne Goraj

Editors: Lisa Duran, Kim Wimpsett, Maureen Adams, Shelby Zimmerman, Alison Moncrieff, Steve Gilmartin, Laura Arendal, Krista Reid-McLaughlin

Compilation Technical Editor: Kirky Ringer

Technical Editors: Matthew Fielder, Kirky Ringer, John Zukowski

Book Designer: Maureen Forsy, Happenstance Type-O-Rama

Graphic Illustrators: Tony Jonick, Patrick Dintino, Inbar Berman

Electronic Publishing Specialists: Cyndy Johnsen and Maureen Forsy

Production Coordinator: Susan Berge

Indexer: Nancy Guenther

Cover Designer: Design Site

Cover Illustrator: Jack D. Myers

SYBEX is a registered trademark of SYBEX Inc.

Mastering, Developer's Handbook, and In Record Time are trademarks of SYBEX Inc.

Screen reproductions produced with Collage Complete.

Collage Complete is a trademark of Inner Media Inc.

Copyright ©1999 SYBEX Inc., 1151 Marina Village Parkway, Alameda, CA 94501. World rights reserved. No part of this publication may be stored in a retrieval system, transmitted, or reproduced in any way, including but not limited to photocopy, photograph, magnetic or other record, without the prior agreement and written permission of the publisher.

Library of Congress Card Number: 99-60006

ISBN: 0-7821-2468-2

Manufactured in the United States of America

10987654

TRADEMARKS:

SYBEX has other descriptive terms

The author and is based upon file based upon pre-publisher make or accuracy of title to performance, of any kind cause

Photographs and file archives and variety of graphical Internet may be of text and image holders, although 17 U.S.C. Section

Netscape Communicator are trademarks

Netscape Communicator this publication is a trademark of Netscape Communications Corporation. All other trademarks

CONTENTS AT A GLANCE

| | |
|--|------------|
| <i>Introduction</i> | xix |
| Part I Introducing Java | 1 |
| Chapter 1 Building the First Java Examples from <i>Java 2: In Record Time</i> by Steve Holzner | 3 |
| Chapter 2 Handling Java Text Fields from <i>Java 2: In Record Time</i> by Steve Holzner | 31 |
| Chapter 3 Using Java Buttons from <i>Java 2: In Record Time</i> by Steve Holzner | 47 |
| Chapter 4 Using Java Layouts and Check Boxes from <i>Java 2: In Record Time</i> by Steve Holzner | 79 |
| Chapter 5 Working with Radio Buttons from <i>Java 2: In Record Time</i> by Steve Holzner | 115 |
| Chapter 6 Adding Scroll Bars from <i>Java 2: In Record Time</i> by Steve Holzner | 149 |
| Part II Java Fundamentals | 177 |
| Chapter 7 Applets, Applications, and the Java Development Kit from <i>Mastering Java 2</i> by John Zukowski | 179 |
| Chapter 8 Working with Java Objects from <i>Mastering Java 2</i> by John Zukowski | 215 |
| Chapter 9 Exception Handling from <i>Mastering Java 2</i> by John Zukowski | 241 |
| Chapter 10 Standard Java Packages from <i>Mastering Java 2</i> by John Zukowski | 265 |
| Chapter 11 File I/O and Streams from <i>Java 2 Developer's Handbook</i> by Simon Roberts and Philip Heller | 311 |
| Part III Advanced Java | 347 |
| Chapter 12 Custom Components from <i>Java 2 Developer's Handbook</i> by Simon Roberts and Philip Heller | 349 |

Chapt

Chapt

Chapt

Chapte

Part IV

Chapter

Chapter

Appendix

Glossary of

Contents at a Glance ix

| | | |
|--------------------------|--|------------|
| Chapter 13 | The JFC Swing Components from <i>Java 2 Developer's Handbook</i> by Simon Roberts and Philip Heller | 409 |
| Chapter 14 | Threads and Multithreading from <i>Mastering Java 2</i> by John Zukowski | 441 |
| Chapter 15 | Java Database Connectivity (JDBC) from <i>Mastering Java 2</i> by John Zukowski | 491 |
| Chapter 16 | The 2D API from <i>Java 2 Developer's Handbook</i> by Simon Roberts and Philip Heller | 537 |
| Part IV | JavaBeans | 593 |
| Chapter 17 | JavaBeans: An Overview from <i>Mastering JavaBeans</i> by Laurence Vanhelsuwe | 595 |
| Chapter 18 | Bean Properties from <i>Mastering JavaBeans</i> by Laurence Vanhelsuwe | 621 |
| Appendix | | 679 |
| | The Essential Java 2 API Reference by David Wall | 681 |
| Glossary of Terms | | 942 |
| Index | | 964 |

TABLE OF CONTENTS

| | |
|---|------------|
| <i>Introduction</i> | <i>xix</i> |
| Part I ► Introducing Java | 1 |
| Chapter 1 ▢ Building the First Java Examples | 3 |
| Building the Hello Example | 4 |
| What's an Applet? | 4 |
| Creating the Hello Example | 5 |
| Setting Up the Java JDK | 6 |
| What's New in 2? | 7 |
| Compiling the Hello Applet | 12 |
| Understanding Java | 13 |
| Running the Hello Applet | 14 |
| Understanding the Hello Example | 15 |
| Object-Oriented Programming | 16 |
| Understanding Java Objects | 16 |
| What's a Java Class? | 17 |
| Learning about Java Packages | 18 |
| Understanding Java Inheritance | 19 |
| What Are Java Access Modifiers? | 21 |
| Understanding the Applet's Web Page | 23 |
| Connecting Java and HTML | 24 |
| What's Next? | 28 |
| Chapter 2 ▢ Handling Java Text Fields | 31 |
| Declaring a Text Field | 33 |
| Initializing with the <i>init()</i> Method | 37 |
| Handling Memory with the <i>new</i> Operator | 38 |
| What Are Java Constructors? | 39 |
| Overloading Java Methods | 40 |
| What's Next? | 44 |

Chapter

Chapter

Chapter 5

| | |
|---|------------|
| Chapter 3 ■ Using Java Buttons | 47 |
| Working with Buttons in Java | 48 |
| Adding a Button to a Program | 50 |
| What Are Java Events? | 52 |
| The <i>this</i> Keyword | 54 |
| Using Button Events | 54 |
| How to Handle Multiple Buttons | 62 |
| Creating <i>clickers.java</i> | 63 |
| Making <i>clickers.java</i> Work | 65 |
| Handling Java Text Areas | 70 |
| Creating <i>txtarea.java</i> | 71 |
| Making <i>txtarea.java</i> Work | 74 |
| What's Next? | 77 |
| Chapter 4 ■ Using Java Layouts and Check Boxes | 79 |
| What Is a Java Layout? | 80 |
| Building the Adder Applet | 80 |
| The <i>Label</i> Control | 82 |
| Adding a Java <i>Label</i> Control | 84 |
| Writing the Adder Applet | 86 |
| Reading Numeric Data from Text Fields | 88 |
| Putting Numeric Data into Text Fields | 90 |
| Working with the Java Grid Layout | 94 |
| Using the <i>GridLayout</i> Manager | 94 |
| Adding a <i>GridLayout</i> Manager | 97 |
| Building Programs with Check Boxes | 101 |
| What's Next? | 113 |
| Chapter 5 ■ Working with Radio Buttons | 115 |
| Building Programs with Radio Buttons | 116 |
| The <i>Radios</i> Applet | 116 |
| Connecting Check Boxes to a <i>CheckboxGroup</i> | 119 |
| Building Programs with Panels | 126 |
| Creating a Panel | 127 |
| Putting Check Boxes and Radio Buttons Together | 132 |
| Creating the Menu Panel | 135 |
| Creating the Ingredients Panel | 136 |

xix

1

3

4

4

5

6

7

12

13

14

15

16

16

17

18

19

21

23

24

28

31

33

37

38

39

40

44

| | |
|--|------------|
| Adding Panels to the <i>sandwich</i> Class | 137 |
| Connecting the Buttons in Code | 139 |
| What's Next? | 147 |
| Chapter 6 ▢ Adding Scroll Bars | 149 |
| Adding Scroll Bars to Programs | 150 |
| Installing Scroll Bars | 152 |
| Connecting Scroll Bars to Code | 154 |
| Using Scroll Bars and BorderLayout | 161 |
| Working with the <i>ScrollPane</i> Class | 172 |
| What's Next? | 176 |
| Part II ► Java Fundamentals | 177 |
| Chapter 7 ▢ Applets, Applications, and the Java Development Kit | 179 |
| Java Applets versus Java Applications | 181 |
| Using the Java Development Kit (JDK) | 183 |
| JDK Utilities | 185 |
| Downloading and Installing the JDK | 187 |
| Building Applications with the JDK | 190 |
| Java Application Source Code | 190 |
| Building Applets with the JDK | 199 |
| HTML for Java Applets | 200 |
| Delivering Applications with the Java Runtime Environment (JRE) | 211 |
| What's New in JDK 1.2 | 212 |
| What's Next? | 213 |
| Chapter 8 ▢ Working with Java Objects | 215 |
| An Introduction to OOP | 216 |
| Data Structures | 216 |
| From Structures to Classes: Encapsulation | 220 |
| Polymorphism | 229 |
| Constructors and Finalizers | 234 |
| Constructors | 234 |
| Garbage Collection | 236 |
| Finalizers | 238 |
| What's Next? | 239 |

Chapter 9

Ov

Exc

Cre

An

Wh

Chapter 10

Jav

Pac

Pac

Pack

137
139
147

149
150
152
154
161
172
176

177

| | |
|---|----------------|
| Chapter 9 ■ Exception Handling | 241 |
| Overview of Exception Handling | 242 |
| The Basic Model | 242 |
| Why Use Exception Handling? | 245 |
| Hierarchy of Exception Classes | 248 |
| Exception-Handling Constructs | 250 |
| Methods Available to Exceptions | 256 |
| The <i>throw</i> Statement | 257 |
| The <i>throws</i> Clause | 257 |
| Creating Your Own Exception Classes | 258 |
| An Example: Age Exceptions | 259 |
| What's Next? | 263 |
| Chapter 10 ■ Standard Java Packages | 265 |
| Java Packages and the Class Hierarchy | 266 |
| Package <i>java.lang</i> —Main Language Support | 268 |
| The Type Wrapper Classes | 269 |
| The String Classes | 271 |
| The <i>Math</i> Library Class | 271 |
| The Multithreading Support Classes | 272 |
| The Low-Level System-Access Classes | 273 |
| The Error and Exception Classes | 274 |
| Package <i>java.util</i> —Utilitarian Language Support | 275 |
| The Core Collection Interfaces | 276 |
| The Concrete Collection Implementation Classes | 277 |
| The Abstract Collection Implementations | 277 |
| The Infrastructure Interfaces and Classes | 278 |
| The Date and Support Classes | 279 |
| The Locale and Supporting Classes | 279 |
| The <i>BitSet</i> Class | 280 |
| The <i>Observer</i> Interface and <i>Observable</i> Class | 280 |
| Package <i>java.io</i> —File and Stream I/O Services | 281 |
| The <i>InputStream</i> Class | 282 |
| The <i>OutputStream</i> Class | 284 |
| The <i>Reader</i> and <i>Writer</i> Classes | 284 |
| The <i>RandomAccessFile</i> Class | 285 |
| The <i>StreamTokenizer</i> Class | 285 |

a

179

181
183
185
187
190
190
199
200
211
212
213

215

216
216
220
229
234
234
236
238
239

Environment (JRE)

| | |
|---|------------|
| Package <i>java.awt</i> —Heart of the Hierarchy | 286 |
| GUI Classes | 288 |
| The Graphics Classes | 293 |
| Geometry Classes | 296 |
| Miscellaneous AWT Classes | 297 |
| Package <i>javax.swing</i> | 298 |
| JComponent Classes | 300 |
| Layout Manager Classes | 302 |
| Model Classes and Interfaces | 303 |
| Manager Classes | 303 |
| <i>AbstractAction</i> and <i>KeyStroke</i> Classes and Action Interface | 303 |
| Miscellaneous Swing Classes | 304 |
| Package <i>java.net</i> —Internet, Web, and HTML Support | 304 |
| Internet Addressing (Classes <i>InetAddress</i> and <i>URL</i>) | 305 |
| Package <i>java.applet</i> —HTML Embedded Applets | 306 |
| Miscellaneous Java Packages | 307 |
| What's Next? | 309 |
| Chapter 11 ■ File I/O and Streams | 311 |
| An Overview of Streams | 312 |
| The Abstract Superclasses | 314 |
| The <i>InputStream</i> Class | 314 |
| The <i>OutputStream</i> Class | 316 |
| The <i>Reader</i> Class | 317 |
| The <i>Writer</i> Class | 318 |
| The Low-Level Stream Classes | 319 |
| The <i>FileInputStream</i> Class | 319 |
| The <i>FileOutputStream</i> Class | 320 |
| The <i>FileReader</i> Class | 321 |
| The <i>FileWriter</i> Class | 321 |
| Other Low-Level Stream Classes | 322 |
| The High-Level Stream Classes | 326 |
| The <i>BufferedInputStream</i> and <i>BufferedOutput</i> <i>Stream</i> Classes | 327 |
| The <i>DataInputStream</i> and <i>DataOutputStream</i> Classes | 328 |
| The <i>LineNumberReader</i> Class | 333 |
| The <i>PrintStream</i> and <i>PrintWriter</i> Classes | 334 |
| The <i>Pushback</i> Classes | 336 |

Part III

Chapter 12

The E

L

E

Strat

C

A

S

D

Subcl

P

P

T

Aggre

T

T

T

Subcl

I

I

T

E

V

What

Chapter 13

A San

J

J

T

286
288
293
296
297
298
300
302
303
303
tion Interface 303
304
ort 304
.URL) 305
306
307
309

311
312
314
314
316
317
318
319
319
320
321
321
322
326

327
lasses 328
333
334
336

| | |
|--|-----|
| The <i>SequenceInputStream</i> Class | 339 |
| The <i>InputStreamReader</i> and <i>OutputStreamWriter</i> Classes | 339 |
| The Non-Stream Classes | 341 |
| The <i>RandomAccessFile</i> Class | 341 |
| The <i>StreamTokenizer</i> Class | 343 |
| What's Next? | 345 |

Part III ► Advanced Java 347

Chapter 12 ◻ Custom Components 349

| | |
|---|-----|
| The Event Delegation Model | 350 |
| Listener Interfaces and Methods | 351 |
| Explicit Event Enabling | 353 |
| Strategies for Designing Custom Components | 354 |
| Component Class Subclassing | 355 |
| Aggregation | 355 |
| Standard Component Subclassing | 356 |
| Design Considerations | 356 |
| Subclassing Component: The <i>Polar</i> Component | 357 |
| <i>Polar</i> 's Look-and-Feel Issues | 357 |
| <i>Polar</i> 's Design Issues | 359 |
| The <i>Polar</i> Component Class and Test Applet | 369 |
| Aggregation: The <i>ThreeWay</i> Component | 376 |
| <i>ThreeWay</i> 's Look-and-Feel Issues | 377 |
| <i>ThreeWay</i> 's Design Issues | 378 |
| The <i>ThreeWay</i> Component and Test Applet | 387 |
| Subclassing a Standard Component: Validating Textfields | 395 |
| <i>IntTextField</i> 's Look-and-Feel Issues | 396 |
| <i>IntTextField</i> 's Design Issues | 396 |
| The <i>IntTextField</i> Component | 399 |
| External Validation: The <i>ValidatingTextField</i> Component | 402 |
| Validating Textfields: Test Applet and Validator Classes | 403 |
| What's Next? | 407 |

Chapter 13 ◻ The JFC Swing Components 409

| | |
|-----------------------------------|-----|
| A Sampler of Swing Components | 410 |
| JFC Textfields, Frames, and Menus | 411 |
| JFC Tabbed Panes | 413 |
| The <i>SwingDemo</i> Class | 415 |

| | | |
|---|------------|------------|
| Improved Components | 418 | JDB |
| JFC Labels | 418 | Cur: |
| JFC Buttons | 419 | Alte |
| JFC Toggles and Check Boxes | 422 | |
| New Components | 424 | |
| JFC Combo Boxes | 424 | |
| JFC Sliders | 426 | |
| JFC Password Fields | 428 | Whz |
| JFC Toolbars | 429 | |
| The <i>SwingDemo</i> Program | 430 | Chapter 16 |
| What's Next? | 439 | The |
| Chapter 14 ■ Threads and Multithreading | 441 | |
| Overview of Multithreading | 442 | Draw |
| Thread Basics | 446 | |
| Creating and Running a Thread | 446 | |
| The Thread-Control Methods | 448 | |
| The Thread Life Cycle | 452 | |
| Thread Groups | 454 | Gen |
| Getting Information about Threads and Thread Groups | 455 | |
| Advanced Multithreading | 459 | |
| Thread Synchronization | 459 | |
| Inter-Thread Communications | 471 | |
| Priorities and Scheduling | 479 | |
| Thread Local Variables | 485 | Whz |
| Daemon Threads | 487 | |
| What's Next? | 488 | Part IV |
| Chapter 15 ■ Java Database Connectivity (JDBC) | 491 | Chapter 17 |
| Java as a Database Front End | 492 | Intri |
| Database Client/Server Methodology | 493 | |
| Two-Tier Database Design | 494 | The |
| Three-Tier Database Design | 495 | |
| The JDBC API | 497 | |
| The API Components | 499 | |
| Limitations Using JDBC (Applications vs. Applets) | 518 | Bea |
| Security Considerations | 520 | |
| A JDBC Database Example | 520 | |
| JDBC Drivers | 528 | |

418
418
419
422
424
424
426
428
429
430
439

441

442

446

446

448

452

454

455

459

459

471

479

485

487

488

491

492

493

494

495

497

499

518

520

520

528

| | |
|---|-----|
| JDBC-ODBC Bridge | 531 |
| Current JDBC Drivers | 531 |
| Alternative Connectivity Strategies | 531 |
| Remote Method Invocation (RMI) | 532 |
| The Common Object Request Broker Architecture (CORBA) | 532 |
| Connectivity to Object Databases | 533 |
| Connectivity with Web-Based Database Systems | 534 |
| What's Next? | 534 |

Chapter 16 ▢ The 2D API 537

| | |
|---|-----|
| The <i>Graphics2D</i> Class and Shapes | 539 |
| The <i>Graphics2D</i> Class | 539 |
| The <i>Shape</i> Interface and Its Implementors | 540 |
| Drawing Operations | 545 |
| Stroking | 545 |
| Filling | 550 |
| Clipping | 556 |
| Transforming | 557 |
| General Paths for Your Own Curves | 561 |
| Specifying a Shape | 561 |
| Transforming a General Path | 562 |
| Drawing a Bezier Curve | 566 |
| Drawing Fractals | 570 |
| Extending the Triangular Fractal | 587 |
| What's Next? | 592 |

Part IV ▸ JavaBeans 593

| | |
|--|-----|
| Chapter 17 ▢ JavaBeans: An Overview | 595 |
| Introduction | 596 |
| What Does a Bean Boil Down to in Practice? | 597 |
| The Black Box View of a Java Bean | 597 |
| Bean Methods | 599 |
| Bean Properties | 599 |
| Bean Events | 600 |
| Bean Environments | 600 |
| Design-Time Environment | 601 |
| Runtime Environment | 603 |
| Applet versus Application Environments | 604 |

| | |
|---|------------|
| The Bean Development Kit and the BeanBox Bean Testing Application | 604 |
| The BeanBox | 605 |
| The BDK Demonstration Beans | 606 |
| Package <i>java.beans</i> | 615 |
| Class <i>Beans</i> | 616 |
| What's Next? | 618 |
| Chapter 18 • Bean Properties | 621 |
| Introduction | 622 |
| The <i>setXXX()</i> and <i>getXXX()</i> Accessor Methods | 623 |
| Defining Read Properties | 624 |
| Defining Write Properties | 625 |
| Defining Read/Write Properties | 625 |
| Bean Property Categories | 627 |
| Simple Properties | 628 |
| Boolean Properties | 628 |
| Indexed Properties | 629 |
| Bound Properties | 630 |
| Constrained Properties | 650 |
| Properties and Multithreading | 665 |
| Multithreading Issues with Simple Properties | 666 |
| Property Listeners and Multithreading | 672 |
| Appendix | 679 |
| The Essential Java 2 API Reference | 681 |
| Glossary of Terms | 942 |
| Index | 964 |

INTRO

Java 2 (the bxx compilati sive cover This book goals in n

► Offe able

► Hel; Java your

► Acq styl thei you

Java 2 you'll nee with Java depths in

If you are many and hard piled rep authors the exha styles. As which ap in comm

You'll authors authors experien evolution

Chapter 1

BUILDING THE FIRST JAVA EXAMPLES

Welcome to Java 2! An ambitious agenda lies before you: You're going to get a firm grip on Java programming, creating both powerful Java programs and Web pages, and you will take a guided tour through Java 2. There is no more exciting programming package available. As you are probably aware, the popularity of Java has skyrocketed as more and more people have seen how versatile and powerful it is. Web programmers have found it an excellent tool because it allows them to write programs that will run on many different types of computers. They have started using it to make their Web pages actually do something.

Java 2

Adapted from *Java 2: In Record Time* by Steven Holzner
ISBN 0-7821-2171-3 560 pages \$29.99

With Java, you will be able to display animation and images, accept mouse clicks and text, use controls like scrollbars and check boxes, print graphics, support pop-up menus, and even support additional windows and menu bars.

We'll start working on your Java skills right away—you won't need to wade through chapters of abstractions first. We will concentrate on examples, on seeing things from the programmer's point of view—on seeing Java at *work*.

Java programs come two ways: as stand-alone applications and as small programs you can embed in Web pages, called *applets*. Of the two, applets are the most popular, and we'll concentrate primarily on them.

BUILDING THE HELLO EXAMPLE

The first example will be a simple one because right now we just want to get you started in Java without too many extra details to weigh you down. You will create a small Java applet, the type of Java program you can embed in a Web page, that will display the words "Hello from Java!"

What's an Applet?

Just what do I mean by an applet? An applet is a special program that you can embed in a Web page such that the applet gains control over a certain part of the Web page. On that part of the page, the applet can display buttons, list boxes, images, and more. Applets make Web pages "come alive."

Each applet is given the amount of space (usually measured in pixels) that it requests in a Web page, such as the amount of space shown in Figure 1.1. (Soon I'll show you how an applet "requests" space.) This is the space that the applet will use for its display. We'll place the words "Hello from Java!" in the applet, as shown in Figure 1.2.

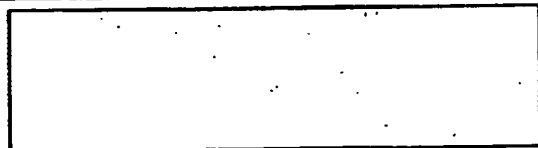


FIGURE 1.1: An applet requests space in a Web page.

FIGURE

That
embed

Creating

Let's ca
lines of
named l
file (suc
files thr
Also not
Word, y
you can
process
to see h
thing bu
hello.

impo
publ

pu

This
see
to d



NC
Not
orth

Hello from Java!

FIGURE 1.2: Hello from Java!

That's how this applet will work; after you create it, you will be able to embed it in a Web page. Let's create and run the applet now.

Creating the Hello Example

Let's call this first applet *hello*. You will store the actual Java code (the lines of text that make up the program) for this applet in a text file named `hello.java`. You'll need an editor of some kind to create this file (such as Windows WordPad or Notepad). You will be creating `.java` files throughout the book, so use an editor you are comfortable with. Also notice that, if you are going to use a word processor like Microsoft Word, you'll have to save your `.java` files as straight text—something you can type out at the DOS prompt and read directly. Check your word processor's Save As menu item or your word processor's documentation to see how to do this. The Sun Java system won't be able to handle anything but straight text files. Now, type the following text into the file `hello.java` (this is the traditional first program in most Java books):

```
import java.awt.Graphics;
public class hello extends java.applet.Applet
{
    public void paint(Graphics g)
    {
        g.drawString("Hello from Java!", 60, 30 );
    }
}
```

This is the text of your first Java program, and soon you'll see what each line means. Having typed in the text, save it to disk as `hello.java`.



NOTE

Note that case counts here—make sure you type `hello.java`, not `Hello.java` or `hello.Java`.

In general, the name of the file will match exactly (including case) the name given in the "class" statement in the file; in this case, that is hello:

```
import java.awt.Graphics;

public class hello extends java.applet.Applet
{
    public void paint( Graphics g )
    {
        g.drawString("Hello from Java!", 60, 30 );
    }
}
```

In this book, you will place your programs into subdirectories of a new directory called java1-2 (this is optional—you can choose any name). That means you'll save the hello.java file as c:\java1-2\hello\hello.java.

Now you have created hello.java. This is the source code for your applet, and it contains the Java code that you have written. The next step is to compile this Java code into a working applet and see your applet at work. Applets have the extension .class, making the name of your actual applet hello.class. I'll show you *why* applets have the extension .class shortly.

SETTING UP THE JAVA JDK

Now you'll use Java itself to create your applet, hello.class, from the code hello.java. If you haven't already done so, you should install the Java Development Kit (JDK) 1.2.

With previous versions of Java, you used to have to go through a rather lengthy and involved installation process, but that's all changed now—you just have to run an .EXE file. You get this .EXE file online, from <http://java.sun.com/products/jdk/1.2/>—just download it and follow the instructions for installation.

The next step is to make sure you can run the JDK from any location in your computer (including the c:\java1-2 directory and its subdirectories, which is where you'll put your Java programs). To do that, make sure the PATH statement in your AUTOEXEC.BAT file (found in the main directory of the C: drive) includes the JDK BIN and LIB directories (here I have installed the JDK in c:\jdk12—use whatever path is appropriate to the way you have installed the JDK):

```
PATH=C:\WINDOWS;C:\JDK12\BIN;C:\JDK12\LIB
```

For W
> Co
Envir
The



You
tory—f
ing a d
long fil



Now
version

What's

If you'r
there to
overview
with Ja
this ma
gramm

From
Many n
by look

Al
pc
de
m
th

For Windows NT, the path will need to be entered into Start > Settings > Control Panel > System. In the System Properties window, select the Environment tab and set the PATH variable.

The JDK 1.2 is ready to go.



TIP

If you need more help installing the JDK, check out the Troubleshooting Web page at <http://www.javasoft.com>.

You can copy the Java documentation from JavaSoft to the same directory—for example, c:\JDK12. Unzip the documentation .zip file, creating a docs subdirectory (your unzipping program must be able to handle long filenames).



NOTE

You'll need a Web browser to look at the Java documentation because it's formatted in HTML.

Now that you've installed Java 2, let's take a look at what's new in this version of Java.

What's New in 2?

If you're familiar with Java 1.0 or Java 1.1, then you'd probably expect there to be some changes in Java 2, and you'd be right. Let's get an overview of the changes in this new edition of Java. If you're not familiar with Java, you should probably skip to the next section and take a look at this material later—much of this won't make any sense unless you've programmed in Java before.

From Java 1.0 to Java 1.1

Many readers will be familiar with Java 1.0, not Java 1.1, so we will start by looking at the changes from Java 1.0 to Java 1.1.

Abstract Windowing Toolkit enhancements Java 1.1 supports printing, faster scrolling, pop-up menus, the clipboard, a delegation-based event model, imaging and graphics enhancements, and more. In addition, it's faster than Java 1.0 (something Java programmers can definitely appreciate)!

.jar files .jar (Java Archive) files were introduced in Java 1.1 and let you package a number of files together, zipping them to shrink them, so the user can download many files at once. You can put many applets and the data they need together into one .jar file, making downloading much faster. These files are analogous to .zip files except that your browser will download them and unzip them on-the-fly for you.

Internationalization Java 1.1 lets you develop *locale-specific applets*, including using Unicode characters, a locale mechanism, localized message support, locale-sensitive date, time, time zone, number handling, and more.

Signed applets and digital signatures Java 1.1 can create digitally signed Java applications. A digital signature gives your users a "path" back to you in case something goes wrong. This is one of the new security precautions popular on the World Wide Web.

Remote method invocation In Java 1.1, RMI lets Java objects have their methods invoked from Java code running in other Java sessions. This is sort of similar to Local Remote Procedure Calls (LRPCs).

Object serialization Serialization was new in Java 1.1, and it lets you store objects and handle them with binary input/output streams. Besides allowing you to store copies of the objects you serialize, serialization is also the basis of communication between objects engaged in RMI. Object serialization is similar to MPC Serialization, for those who are familiar with Microsoft's Foundation Classes.

Reflection In Java 1.1, reflection lets Java code examine information about the methods and constructors of loaded classes and make use of those reflected methods and constructors.

Inner classes Java 1.1 makes it easier to create adapter classes. An adapter class is a class that implements an interface required by an API (Applications Programming Interface). An adapter class "delegates" control back to an enclosing main object.

New Java native method interface Native code is code that is written specifically for a particular machine. In Java 1.1, this interface was introduced to provide a standard programming

That
idea w

From
Now l

interface for writing Java native methods. The primary goal is binary compatibility of native method libraries across all Java virtual machine implementations on a given platform. Writing and calling native code can significantly improve execution speeds. Java 1.1 included a powerful new Java native method interface.

Byte, Short, and Void classes In Java 1.1, Byte and Short values can be handled as "wrapped" numbers when you use the new Java classes Byte and Short. The new Void class is a placeholder class that we can derive classes from, rather than use directly.

Deprecated methods Quite a number of Java 1.0 methods were considered obsolete in Java 1.1, and they are marked as deprecated in the Java 1.1 documentation. (The Java compiler now displays a warning when it compiles code that uses a deprecated feature.)

Networking enhancements Networking enhancements in Java 1.1 included support for selected BSD-style socket options in the `java.net` base classes. With Java 1.1, `Socket` and `ServerSocket` are non-final, extendable classes. New subclasses of `SocketException` were added for finer granularity in reporting and handling network errors.

I/O enhancements In Java 1.1, the I/O package was extended with character streams, which are like byte streams except that they contain 16-bit Unicode characters rather than eight-bit bytes. Character streams make it easy to write programs that are independent of a specific character encoding and are therefore easier to internationalize. Nearly all of the functionality available for byte streams is also available for character streams.

That completes this overview of what's new in Java 1.1—if you have no idea what I'm talking about, don't worry, it'll become clear later.

From Java 1.1 to Java 2

Now let's have a look at what's new in Java 2.

Security enhancements When code is loaded, it is assigned permissions based on the security policy currently in effect. Each permission specifies a permitted access to a particular

resource (such as "read" and "write" access to a specified file or directory, "connect" access to a given host and port, and so on). The policy, specifying which permissions are available for code from various signers/locations, can be initialized from an external configurable policy file. Unless a permission is explicitly granted to code, it cannot access the resource that is guarded by that permission.

Swing (JFC) Swing is the part of the Java Foundation Classes (JFC) that implements a new set of GUI components with a "pluggable" look and feel. Swing is implemented in pure Java, and is based on the JDK 1.1 Light-weight UI Framework. The pluggable look and feel lets you design a single set of GUI components that can automatically have the look and feel of any platform (e.g., Windows, Solaris, Macintosh).

Java 2D (JFC) The Java 2D API is a set of classes for advanced 2D graphics and imaging. It encompasses line art, text, and images in a single comprehensive model.

Accessibility (JFC) Through the Java Accessibility API, developers will be able to create Java applications that can interact with assistive technologies such as screen readers, speech recognition systems, and Braille terminals.

Drag and Drop (JFC) Drag and Drop enables data transfer across both Java and native applications, between Java applications, and within a single Java application.

Collections The Java Collections API is a unified framework for representing and manipulating Java collections (I'll show you more about them later), allowing them to be manipulated independent of the details of their representation.

Java extensions Framework Extensions are packages of Java classes (and any associated native code) that application developers can use to extend the core platform. The extension mechanism allows the Java Virtual Machine (JVM) to use the extension classes in much the same way it uses the system classes.

JavaBeans enhancements Java 2 provides developers with standard means to create more sophisticated JavaBeans components and applications that offer their customers more seamless integration with the rest of their runtime environment,

suc
bro

Imp
enal
nese

Paci
pac
can
Envi

RM
seve
whic
obje
rem
use
such

Seri
API
inde
data
exist
writi

Ref
ence
exam
obje
that
Java

Audi
soun
appl

Java
Brok
base
tribu
trans

pecified file or
rt, and so on).
lable for code
from an exter-
explicitly
t is guarded

dation Classes
nts with a
in pure Java,
ework. The
of GUI com-
feel of any

es for
as line art,

lity API,
that can
readers,

its transfer
Java applica-

l framework
(I'll show
manipulated

ages of Java
rtion devel-
ision mecha-
he extension
ses.

opers with
eans compo-
ore seam-
nment,

such as the desktop of the underlying operating system or the browser.

Input method framework The input method framework enables all text-editing components to receive Japanese, Chinese, or Korean text input through standard input methods.

Package version identification "Versioning" introduces package level version control where applications and applets can identify (at runtime) the version of a specific Java Runtime Environment, VM, and class package.

RMI enhancements Remote Method Invocation (RMI) has several new enhancements including Remote Object Activation, which introduces support for remote objects and automatic object activation, as well as Custom Socket Types that allow a remote object to specify the custom socket type that RMI will use for remote calls to that object. (RMI over a secure transport, such as SSL, can be supported using custom socket types.)

Serialization enhancements Serialization now includes an API that allows the serialized data of an object to be specified independently of the fields of the class. This allows serialized data fields to be written to and read from the stream using the existing techniques (this ensures compatibility with the default writing and reading mechanisms).

Reference objects A reference object encapsulates a reference to some other object so that the reference itself may be examined and manipulated like any other object. Reference objects allow a program to maintain a reference to an object that does not prevent the object from being reclaimed by the Java "garbage collector," which manages memory.

Audio enhancements Audio enhancements include a new sound engine and support for audio in applications as well as applets.

Java IDL Java IDL adds CORBA (Common Object Request Broker Architecture) capability to Java, providing standards-based interoperability and connectivity. Java IDL enables distributed Web-enabled Java applications to invoke operations transparently on remote network services using the industry

standard OMG IDL (Object Management Group Interface Definition Language) and IIOP (Internet Inter-ORB Protocol) defined by the Object Management Group.

JAR enhancements The enhancements include added functionality for the command-line JAR tool for creating and updating signed JAR files. There are also new standard APIs for reading and writing JAR files.

JNI enhancements The Java Native Interface (JNI) is a standard programming interface for writing Java native methods and embedding the Java Virtual Machine into native applications. The primary goal is binary compatibility of native method libraries across all Java Virtual Machine implementations on a given platform. Java 2 extends the Java Native Interface to incorporate new features in the Java platform.

JVMDI A new debugger interface, the Java Virtual Machine, now provides low-level services for debugging. The interface for these services is the Java Virtual Machine Debugger Interface (JVMDI).

JDBC enhancements Java Database Connectivity (JDBC) is a standard SQL database access interface, providing uniform access to a wide range of relational databases. JDBC also provides a common base on which higher-level tools and interfaces can be built. The Java 2 software bundle includes JDBC and the JDBC-ODBC bridge.

These concepts will become clearer as we proceed. Now, you're ready to compile the hello applet and see it at work.

Compiling the Hello Applet

Now that you have installed the JDK and have your `hello.java` source file ready to go, you can create the actual applet and see it run. To do this, change to the `c:\java1-2\hello` directory now (or wherever you have saved the `hello.java` file); this is how the DOS prompt should look:

```
c:\java1-2\hello>
```

Next, type this to create your applet:

```
c:\java1-2\hello>javac hello.java
```

Th
into .
comp
DOS
see be
hell
What

UNDE

Let's t
ming l
then t
unders
it such
your p
use. In
named
Java-cc
this wa
a few b
browse
that is
those a
so othe
shortly.

Expe
isn't Ja
comput
machin
machin
net—it o
long as
are run
convert
puters c

The J
always s
support

The name of the Java program that takes your Java code and turns it into `.class` files ready to run in Web pages is `javac.exe`, the Java compiler (i.e., it compiles `.java` files into `.class` files). If you type the DOS command `Dir` to look at the current directory contents, you should see both `hello.java` and `hello.class`. Because you've created `hello.class`, your applet is ready to go—but what does that mean? What have you really done?

UNDERSTANDING JAVA

Let's take the time now to get an overview of Java. As in most programming languages, we write Java code using words and numbers that are then translated—that is, *compiled*—into binary files that computers can understand. The `hello.java` program is an example of this—you write it such that you can understand it, but when you want to actually run your program, you have to compile it into something a computer can use. In this case, that means using the Java compiler to produce the file named `hello.class`. `hello.class` is a binary file of *bytecodes* that Java-compatible Web browsers can run to produce the desired result. In this way, several lines of Java program code can be compiled neatly into a few bytes. Those bytes are what is actually downloaded when Web browsers read the Web page in which you have placed your Java applets—that is to say, the actual applet is a `.class` file, like `hello.class`, and those are the files you place on your Internet Service Provider's server so other people's Web browsers can download them, as you'll see very shortly.

Experienced programmers may wonder about these bytecodes—why isn't Java simply compiled into the normal machine code that each computer really runs? Because Java bytecodes were intentionally made machine-independent so that they could be run on a wide variety of machines, and that is what originally made them so popular on the Internet—it doesn't matter what type of machine you're downloading to, as long as the user's Web browser can run Java. The downloaded bytecodes are run by the *Java Virtual Machine*, or JVM, and it is the JVM's task to convert bytecodes into the machine language that users' individual computers can run.

The JVM is actually a hypothetical chip that runs Java—it is almost always software, not hardware, that runs Java. Each Web browser that supports Java has a JVM built right into it, and it loads the `.class` file

erface Defi-
pool)

dded func-
and updat-
ls for

l) is a stan-
methods
applica-
ve method
ions on a
ce to incor-

Machine,
terface for
interface

(JDBC) is
uniform
also pro-
interfaces
IC and the
you're ready

java source
in. To do this,
ever you have
ould look:

that makes up your applet with JVM's *class loader* and then runs the applet.

Running the Hello Applet

To see `hello.class`, your first applet, running, you'll need a Web page to place it in. Use your editor again to create a new file, `hello.htm`, which will be your Web page, written in the language of Web pages, *HyperText Markup Language* (HTML) (we'll review HTML in a minute). Enter the following text into `hello.htm` and save it in the same directory as the `hello.class` file:

```
<html>

<!-- Web page written for the Sun Applet Viewer -->

<head>
<title>hello</title>
</head>

<body>
<hr>

<applet
code=hello.class
width=200
height=200>

</applet>
<hr>
</body>
</html>
```

Now you can run the hello applet by simply viewing this new Web page, `hello.htm`. To do that, use the Applet Viewer that comes with the JDK 1.2. To use the Applet Viewer, go back to the `hello` subdirectory and type the following:

```
c:\java1-2\hello>appletviewer hello.htm
```

Again, capitalization is very important here—make sure your capitalization matches the exact spelling of the Web page name. When you've done

thi
"H.

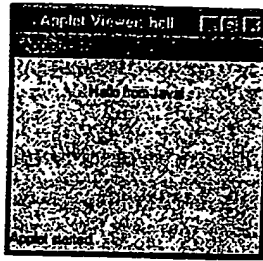


Y
Let
exa
wor.

Under
Let

V
adv
rate
grea
enti
Inst
fin
tine
doin

this, the Applet Viewer runs, as shown below—and you see your message, “Hello from Java!” Your first applet is a success.



TIP

You can use any Java-enabled Web browser to look at this Web page. For most of the applets in this book, however, you will have to use either a Web browser that supports Java 2 (not just Java 1.0 or Java 1.1) or the Sun Applet Viewer.

Your first applet, `hello.class`, runs—but what exactly did you do? Let's take a look now at the Java code that you entered for `hello.java`, examining it line by line to get a better idea of how Java programming works (even though Java will handle many of these details for you later).

Understanding the Hello Example

Let's take apart your first applet now. Begin with this line:

```
import java.awt.Graphics;
```

What does this mean? This line actually points out one of the great advantages of Java programming. When you're adding menus and separate windows to your Java applets, you can imagine that it would be a great deal of work to create everything from scratch—that is, write the entire code for menu handling, separate window creation, and so forth. Instead of asking you to do so, Java comes complete with several predefined libraries, and much of this book will be an examination of the routines in these libraries. You'll learn more about this later, but what you're doing is adding support from the main Java graphics library of routines

to your applet. In this way, we'll be able to draw the text string, "Hello from Java!", in the applet's window.

**NOTE**

If you're a C/C++ programmer, you'll notice that the import statement works much like the C/C++ #include statement.

Next, add these lines to `hello.java`:

```
import java.awt.Graphics;

public class hello extends java.applet.Applet
```

You've just created a Java *class* named `hello`. What does this mean?

OBJECT-ORIENTED PROGRAMMING

Objects and *classes* are two fundamental concepts in object-oriented languages like Java. There's been a lot of hype about object-oriented programming (OOP), and that can make the whole topic seem mysterious and unapproachable. In fact, object-oriented programming was introduced to make longer programs *easier* to create. We'll start a mini-survey of object-oriented programming by looking at objects.

Understanding Java Objects

In long, involved programs, there can be a profusion of both variables and functions, sometimes hundreds of each. Creating and maintaining the program code can become a very cluttered task because you have to keep so many things in mind. There may also be unwanted interaction if various functions use variables of the same name. Object-oriented programming was invented to break up such large programs.

The idea behind objects is quite simple—you just break up your program into the various parts, each of which you can easily conceptualize as performing a discrete task, and those are your objects. For example, you may put all the screen-handling parts of a program together into

an ob
tions
varial
scre
but al
or dr
the re
gram
As
less u
forth
matic
usefu
way i

What's

But h
an ob
term
data
you s

Th
the t
This
you n



For
you c

Ye
grap
just

an object named `screen`. Objects are more powerful than simple functions or sets of variables because an object can hold both functions and variables wrapped up together in a way that makes it easy to use. The `screen` object may hold not only all the data displayed on the screen, but also the functions needed to handle that data, like `drawString()` or `drawLine()`. This means that all the screen handling is hidden from the rest of the program in a convenient way, making the rest of the program easier to handle.

As another example, think of a refrigerator. A refrigerator would be far less useful if you had to regulate all the temperatures and pumps and so forth by hand at all times. Making all those functions internal and automatic to the refrigerator makes it into an easy object to deal with and a useful one: a *refrigerator*. Wrapping up code and data into objects this way is the basis of object-oriented programming.

What's a Java Class?

But how do you create objects? That's where *classes* come in. A class is to an object what a cookie cutter is to a cookie—a template or blueprint. In terms of programming, you might think of the relationship between a data type, like an integer, and the actual variable itself like this, where you set up an integer named `the_data`:

```
int the_data;
```

This is the actual way to create an integer variable in Java. Here, `int` is the type of variable you are declaring and `the_data` is the variable itself. This is the same relationship that a class has to an object, and informally you may think of a class as an object's type.



TIP

Java supports all the standard C and C++ primitive data types like `int`, `double`, `long`, `float`, and so forth.

For example, if you had set up a class named, say, `graphicsClass`, you can create an object of that class named `screen` this way:

```
graphicsClass screen;
```

You'll see how to actually create a class soon (creating a class like `graphicsClass` is not hard—when you create a class in code, you will just group all its functions and data inside the class definition), and then

you'll see how to create objects of that class. What's important to remember is this: the object itself is what holds the data you want to work with; the class itself holds no data but just describes how the object should be set up.

Object-oriented programming at root is nothing more than a way of grouping functions and the data they work on together to make your program less cluttered. You'll see more about object-oriented programming throughout this book, including how to create a class, how to create an object of that class, and how to reach the functions and data in that object when you want to.

That completes the mini-overview of classes and objects. As you can see, a class is just a programming construct that groups together, or *encapsulates*, functions and data, and an object may be thought of as a variable of that class's type, as the object `screen` is to the class `screenClass`.

As it turns out, Java comes complete with several libraries of predefined classes, which save you a great deal of work. Throughout this book, we will examine these predefined and very useful Java classes. Using these predefined classes, we'll create objects needed to handle buttons, text fields, scroll bars, and much more.

Learning about Java Packages

These class libraries are called *packages* in Java, and one such library is called `java.awt` (where `awt` stands for Abstract Window Toolkit). This library holds the `Graphics` class, which will handle the graphics work you undertake. So this line in the `hello.java` file:

```
import java.awt.Graphics;
```

actually means that you want to include the Java Graphics class and make use of it in your program. In a minute, you will use an object of the `Graphics` class for your graphics output.

You've added support for graphics handling by including the `java.awt.Graphics` class (and in Java, displaying the text string "Hello from Java!" is considered graphics handling). Next, it's time to set up your hello applet itself. To do so, define a new class named `hello`. This is the standard way of setting up an applet in Java, and in fact, the applet itself has the file extension `.class`. That's because each class defined in a `.java` file ends up being exported to a `.class` file, where you can make use of it. You'll learn more details about this soon.

Und

y
t
j
n
y
w
c
ir
cl
x
st
b
t
it
t
/
cl
de
us

It would be quite difficult to write all the code an applet class needs from scratch. For example, we'd need to interact with the Web browser, reserve a section of screen, initialize the appropriate Java packages, and much more. It turns out that all that functionality is already built into the Java Applet class, which is part of the `java.applet` package. But how do you make use of the Applet class? You want to customize the applet to display your text string, and the `java.applet.Applet` class itself knows nothing about that.

Understanding Java Inheritance

You can customize the `java.applet.Applet` class by *deriving* the `hello` class from the `java.applet.Applet` class. This makes `java.applet.Applet` the *base* class of the `hello` class, and it makes `hello` a class derived from `java.applet.Applet`. This gives you all the power of the `java.applet.Applet` class without the worries of writing it yourself, and you can add what you want to this class by adding code to your derived class `hello`.

This is an important part of object-oriented programming, and it's called *inheritance*. In this way, a derived class inherits the functionality of its base class and adds more on top of it. For example, you may have a base class called `chassis`. You can derive various classes from this base class called, say, `car` and `truck`. In this way, two derived classes can share the same base class, saving time and effort programmatically. Although the `car` and `truck` classes share the same base class, `chassis`, they added different items to the base class, ending up as two quite different classes, `car` and `truck`.

Using inheritance, then, you will *extend* the base class `java.applet.Applet` by creating your own class `hello` and adding onto the base class. In the `hello.java` source, you indicate that the `hello` class is derived from the `java.applet.Applet` class like this (note that you use the keyword `class` to indicate that you are defining a new class):

```
import java.awt.Graphics;
```

```
public class hello extends java.applet.Applet
```

```
{
```


In starting to set up the new class, `hello`, you've given it all the power of the `java.applet.Applet` class (like the ability to request space from the Web browser and to respond to many browser-created commands). But how do you make additions and even alterations to the `java.applet.Applet` class to customize your own `hello` class? How do you display your text string? One way is by *overriding* the base class's built-in functions (overriding is an important part of object-oriented programming). When you redefine a base class's function in a derived class, the new version of the function is the one that takes over. In this way, you can customize the functions from the base class as you like them in the derived class.

For example, one function in the `java.applet.Applet` class is called `paint()`. This is a very important function that is called when the Web browser tells the applet to create its display on the screen. This happens when the applet first begins and every time it has to be redisplayed later (for example, if the Web browser was minimized and then maximized, or if some window was moved and the applet's display area was uncovered after having been covered).

Your goal in the `hello` class is to display the string "Hello from Java!" on the screen, and in fact, you will override the `java.applet.Applet` class's `paint()` function to do so. You override a base class's function simply by redefining it in the new class. Do that now for the `paint()` function, noting first that the built-in functions of a class are called that class's *methods*. In this case, then, you override (that is, redefine) the `paint()` method like this:

```
import java.awt.Graphics;

public class hello extends java.applet.Applet
{
    public void paint( Graphics g ){
```



NOTE

The built-in functions of a class are called *methods*. Classes can also have built-in variables—called *data members*—and even constants. Collectively, all these parts are called a class's *members*.

What

Th
de
th
in
on
be
de

a n
In
ret
(th
ab
for
ob

cal

Th
In
sc
th

na
wt
str
in
"p
etc
pk

What Are Java Access Modifiers?

The keyword *public* is called an *access modifier*. A class's methods can be declared *public*, *private*, or *protected*. If they are declared *public*, then you can call them from anywhere in the program, not just in the class in which they are defined. If they are *private*, they may be called from only the class in which they are defined. If they are *protected*, they may be called from only the class in which they are defined and the classes derived from that class.

Next, indicate the *return* type of the `paint()` method. When you call a method, you can pass parameters to it, and it can return data to you. In this case, `paint()` has no return value, which you indicate with the return type *void*. Other return types are *int* for an integer return value (this variable is usually 32 bits long), *long* for a long integer (this variable is usually 64 bits long), *float* for a floating point return, or *double* for a double-precision floating point value. You can also return arrays and objects in Java.

Finally, note that you indicate that the `paint()` method is automatically passed one parameter—an object of the *Graphics* class called *g*:

```
import java.awt.Graphics;

public class hello extends java.applet.Applet
{
    public void paint(Graphics g)
    {
```

This *Graphics* object represents the physical display of the applet. That is, you can use the built-in methods of this object—such as `drawImage()`, `drawLine()`, `drawOval()`, and others—to draw on the screen. In this case, you want to place the string "Hello from Java!" on the screen, and you can do that with the `drawString()` method.

How do you reach the methods of an object like the *Graphics* object named *g*? You do that with a dot operator (`.`) like this: `g.drawString()`, where here you are invoking *g*'s `drawString()` method to "draw" a string of text on the screen (text is handled like any other type of graphics in a windows environment—that is, it is drawn on the screen rather than "printed," just as you would draw a rectangle or circle). Supply three parameters to the `drawString()` method—the string of text you want to display, and the (*x*, *y*) location of that string's lower-left corner (called the

starting point of the string's *baseline*) in pixels on the screen, passed in two integer values. As shown in Figure 1.3, you can draw your string at the pixel location (60, 30), where (0, 0) is the upper-left corner of the applet's display.

**NOTE**

The coordinate system in a Java program is set up with the origin (0, 0) at the upper left, with x increasing horizontally to the right and y increasing vertically downwards; this fact will be important throughout the book. If it seems backwards to you, you might try thinking of it in terms of reading a page of text, like this one, where you start at the upper-left and work your way to the right and down. The units of measurement in Java coordinate systems are almost always screen pixels.

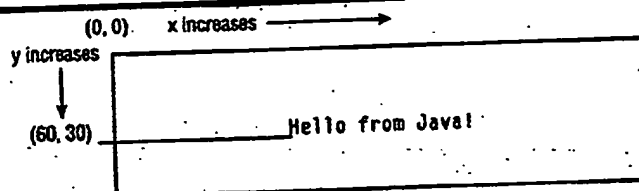


FIGURE 1.3: Drawing a string at (60,30)

This means that you add a call to the `drawString()` method this way:

```
import java.awt.Graphics;
public class hello extends java.applet.Applet
{
    public void paint( Graphics g )
    {
        g.drawString( "Hello from Java!", 60, 30 );
    }
}
```

Note that Java uses the same convention as C or C++ to indicate that a code statement is finished: it ends the statement with a semicolon (;).

**TIP**

In general, Java adheres very strongly to C++ coding conventions. If you know C++, you already know a great deal of Java.

Y
to s
Java
class
your
page
ates
if ap
B
all y
hel
View

Under

The
Web

We

FIGU

He
open
apple

<
<
<

You have completed the code necessary for this applet, which is also to say you have completed the code for the new class, `hello`. When the Java compiler creates `hello.class`, the entire specification of the new class will be in that file. This is the actual binary file that you upload to your Internet Service Provider so that it may be included in your Web page. A Java-enabled Web browser takes this class specification and creates an object of that class and then gives it control to display itself and, if applicable, handle user input.

But how? You have not yet completed the dissection of the first example; all you have done so far is to trace the development of `hello.java` into `hello.class`. How did you get the applet to be displayed in the Applet Viewer?

Understanding the Applet's Web Page

The Applet Viewer took the `hello.class` applet and displayed it in a Web page, as shown in Figure 1.4.

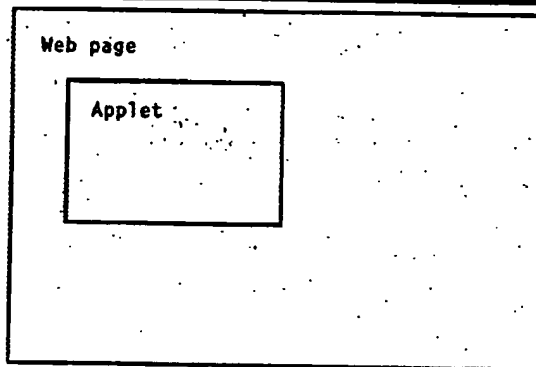


FIGURE 1.4: Displaying an applet in a Web page

How did it get there? You created a Web page for your applet and then opened that Web page in the Applet Viewer, which then displayed your applet. That Web page looks like this:

```
<html>
<!-- Web page written for the Sun Applet Viewer -->
<head>
```

```
<title>hello</title>
</head>
```

```
<body>
<hr>
```

```
<applet
code=hello.class
width=200
height=200>
```

```
</applet>
```

```
<hr>
</body>
</html>
```

Web pages are written in HTML (HyperText Markup Language). Because applets appear in Web pages, we will take the time to briefly work through the above page to make sure you know what's going on. If you're familiar with HTML, you can skip much of this review, but you should take a look at how to use the `<applet>` tag to embed applets in Web pages.

CONNECTING JAVA AND HTML

Let's take apart the Web page you created for the applet now, starting with the `<html>` tag:

```
<html>
```

Instructions in .html pages are placed into tags surrounded by angle brackets: `<` and `>`. The tags hold directions to the Web browser and are not displayed on the screen. Here, the `<html>` tag indicates to the Web browser that this .html file is written in HTML.

Next comes a comment. Comments in .html pages are written using the `!` symbol like this: `<!-- This is a comment.-->`. Indicate that this is a Web page written so that we can use the Sun Applet Viewer, like this:

```
<html>
```

```
<!-- Web page written for the Sun Applet Viewer-->
```

wit
enc
suc
tex
the

1
the
is w
a ru

N
tag,
port
200

Next comes the header portion of the Web page, which you declare with the `<head>` tag, ending the header section with the corresponding end header tag, `</head>` (many HTML tags are used in pairs like this, such as `<head>` and `</head>`, or `<center>` and `</center>` to center text and images). In this case, the `.html` file gets the title (set up with the `<title>` tag) *hello*, to match your applet:

```
<html>

<!-- Web page written for the Sun Applet Viewer -->
```

```
<head>
<title>hello</title>
</head>
```

The title is the name given to a Web page, and it's usually displayed in the Web browser's title bar. Next comes the body of the Web page. Here is where all the actual items for display will go. You start the page off with a ruler line (visible in Figure 1.4), using the `<hr>` tag:

```
<html>

<!-- Web page written for the Sun Applet Viewer -->
```

```
<head>
<title>hello</title>
</head>
```

```
<body>
<hr>
```

Now we come to the applet. Applets are embedded with the `<applet>` tag, and here you use the `code` keyword to indicate that this applet is supported by the `hello.class` file. You indicate the size of the applet as `200 x 200` pixels (you can choose any size you like here) this way:

```
<html>

<!-- Web page written for the Sun Applet Viewer -->
```

Language).
ime to briefly
at's going on. If
view, but you
mbed applets in

now, starting

ounded by angle
rowser and are
stes to the Web

re written using
licate that this is
ewer, like this:

```

<head>
<title>hello</title>
</head>

<body>
<hr>
<applet
code=hello.class
width=200
height=200>
</applet>

```

**TIP**

You can also use the `java.applet.Applet.resize()` method in your source code to request that the Web browser resize applets.

The `<applet>` tag is important, so let's take a closer look at it now. Here's how the `<applet>` tag works in general (the items in square brackets are optional, and the others are required):

```

<APPLET>
  [ALIGN = LEFT or RIGHT or TOP or TEXTTOP or MIDDLE or
  ABSMIDDLE or BASELINE or BOTTOM or ABSBOTTOM]
  [ALT = AlternateText]
  CODE = AppletName.class
  [CODEBASE = URL of .class file]
  HEIGHT = AppletPixelsHeight
  [HSPACE = PixelSpaceToLeftOfApplet]
  [NAME = AppletInstanceName]
  [VSPACE = PixelSpaceAboveApplet]
  WIDTH = AppletPixelsWidth
  >
  [<PARAM NAME = Parameter1 VALUE = VALUE1]
  [<PARAM NAME = Parameter2 VALUE = VALUE2]

```

```

</APPLET>

```



ap
et
tc
al
yc



br
pk
pk

Ja
</

**TIP**

You can specify the URL of the applet's .class file with the CODEBASE keyword. This is often useful if you want to store your applets together in a directory in your ISP, away from the .html files.

Indicate to the Web browser here how much space you'll need for your applet, using the HEIGHT and WIDTH keywords. You can also pass parameters to applets with the PARAM keyword like this: `<applet> PARAM today = "friday" </applet>`. Passing parameters in this way allows you to customize your applets to fit different Web pages because you can read the parameters from inside an applet and make use of them.

**TIP**

There are enhancements to the `<applet>` tag in Java 2, such as the ability to pass the name of .jar files as parameters. You'll learn more about this later on.

Not all Web browsers support Java. In practice, this means that those browsers just ignore the `<applet>` tag. This, in turn, means that you can place text between the `<applet>` and `</applet>` tags that will be displayed in non-Java browsers (and not in Java-enabled browsers), like this:

```
<applet code=hello>
```

```
Your Web browser does not support Java, so you can't see my
applets, sorry!
```

```
</applet>
```

Using the `<applet>` tag, you can embed applets in Web pages, as Java has done in this temporary page. Finish off the Web page with the `</body>` and `</html>` tags as follows:

```
<html>
```

```
<!-- Web page written for the Sun Applet Viewer>
```

```
<head>
```

```
<title>hello</title>
```

```
</head>
```

```
<body>
```

```
<hr>
```

```
<applet
```

```
code=hello.class
```

```
width=200
```

method in your source

look at it now.
ems in square

MIDDLE or
TOM]


```
height=200>
```

```
</applet>
```

```
<hr>
```

```
</body>
```

```
</html>
```

This completes our first example—you've had a glimpse into the process of creating and running an applet. It was as quick and easy as that—you created and ran your first applet.

WHAT'S NEXT?

In this chapter, the example applet demonstrated the easiest way to get an applet to work. Let's continue on to get a better idea of how you'll be working with Java throughout the book as you give your applet more power in Chapter 2.

Java 2 COMPLETE

Java 2 Complete is the most comprehensive book available for both the beginner and the experienced Java programmer. It covers everything you need to know to get started with Java 2, from the basics of the language to the advanced topics of applets, applications, and the Java Development Kit.

Java 2 Complete covers everything you need to know to get started with Java 2, from the basics of the language to the advanced topics of applets, applications, and the Java Development Kit. It includes everything you need to know to get started with Java 2, from the basics of the language to the advanced topics of applets, applications, and the Java Development Kit.

Java 2 Complete is the most comprehensive book available for both the beginner and the experienced Java programmer. It covers everything you need to know to get started with Java 2, from the basics of the language to the advanced topics of applets, applications, and the Java Development Kit.

USER LEVEL ALL LEVELS

BOOK TYPE HOW-TO/REFERENCE

CATEGORY INTERNET PROGRAMMING



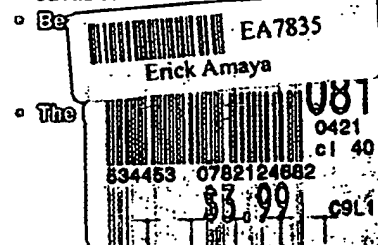
Inside:

- Building the First Java Examples
- Handling Java Text Fields
- Using Java Buttons
- Using Java Layouts and Check Boxes
- Working with Radio Buttons
- Adding Scroll Bars

- Applets, Applications, and the Java Development Kit
- Working With Java Objects
- Exception Handling
- Standard Java Packages
- File I/O and Streams

- Custom Components
- The JFC Swing Components
- Threads and Multithreading
- Java Database Connectivity (JDBC)
- The 2D API

- JavaBeans: An Overview



A New Public Key Cryptosystem Based on Higher Residues

✓ David Naccache

Gemplus Card International

34 rue Guynemer

Issy-les-Moulineaux CEDEX, 92447, France

naccache@compuserve.com

Jacques Stern

Ecole Normale Supérieure

45 rue d'Ulm

Paris CEDEX 5, 75230, France

jacques.stern@ens.fr

Abstract

This paper describes a new public-key cryptosystem based on the hardness of computing higher residues modulo a composite RSA integer. We introduce two versions of our scheme, one deterministic and the other probabilistic. The deterministic version is practically oriented: encryption amounts to a single exponentiation w.r.t. a modulus with at least 768 bits and a 160-bit exponent. Decryption can be suitably optimized so as to become less demanding than a couple RSA decryptions. Although slower than RSA, the new scheme is still reasonably competitive and has several specific applications. The probabilistic version exhibits an homomorphic encryption scheme whose expansion rate is much better than previously proposed such systems. Furthermore, it has semantic security, relative to the hardness of computing higher residues for suitable moduli.

1 Introduction

It is striking to observe that two decades after the discovery of public-key cryptography, the cryptographer's toolbox still contains very few asymmetric encryption schemes. Consequently, the search for new public-key mechanisms remains a major challenge. The quest appears sometimes hopeless as new schemes are immediately broken or, if they survive, are compared with RSA, which is obviously elegant, simple and efficient.

Similar investigations have been relatively successful in the related setting of identification, where a user attempts to convince another entity of his identity by means of an on-line communication. For example, there have been several attempts to build identification protocols based on simple operations (see [33, 35, 36, 26]). Although the question of devising new public-key cryptosystems appears much more difficult (since it deals with trapdoor functions rather than simple one-way functions), we feel that research in this direction is still in order: simple yet efficient constructions may have been overlooked.

The scheme that we propose in the present paper uses an RSA integer n which is a product of two primes p and q ,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

5th Conference on Computer & Communications Security
San Francisco CA USA

Copyright ACM 1998 1-58113-007-4/98/11...\$5.00

as usual. However, it is quite different from RSA in many respects:

1. it encrypts messages by exponentiating them with respect to a fixed base rather than by raising them to a fixed power
2. it uses a different "trapdoor" for decryption
3. its strength is not directly related to the strength of RSA
4. it exhibits further "algebraic" properties that may prove useful in some applications.

We briefly comment on those differences. The first one may offer a competitive advantage in environments where a large amount of memory is available: such environments allow impressive speed-ups in exponentiations that do not have analogous counterparts in RSA-like operations. The second is of obvious interest in view of the fact quoted above that there are very few public-key cryptosystems available. Without going into technical details at this point, let us simply mention that the new trapdoor is obtained by injecting small prime factors in $p-1$ and $q-1$. In order to understand what the third difference is, we note that, if the modulus n can be factored, then both RSA and the proposed cryptosystem are broken. However, it is an open problem whether or not RSA is "equivalent" to factoring, which would mean that breaking RSA allows to factor. For this reason, the hypothesis that RSA is secure has become an assumption of its own, formally stronger than factoring. Our cryptosystem is related to another hypothesis, also formally stronger than factoring and known as the higher residuosity assumption. This may help to understand how these various hypotheses are related. Finally, we will explain the algebraic property of our scheme (called the *homomorphic property*) by means of an example: suppose that one wishes to withdraw a small amount u from the balance m of some account; assume further that the balance is given in encrypted form $E(m)$ and that the clerk performing the operation does not have access to decryption. The cryptosystem that we propose simply solves the problem by computing $E(m)/E(u) \bmod n$, which turns out to be the encryption of the new balance $m - u$.

The ability to perform algebraic operations such as additions or subtractions by playing only with the cryptograms has potential applications in several contexts. We quote a few:

1. in election schemes, it provides a tool to obtain the tally without decrypting the individual votes (see [4])

2. in the area of watermarking, it allows to add a mark to previously encrypted data (as explained in [25]).

Still, in these contexts, it is often needed to encrypt data taken from a small set S (e.g. 0/1 votes) and it is well known that deterministic cryptosystems, such as RSA, fail here: in order to decrypt $E(a)$, one can simply compare the ciphertext with the encryptions of all members of S and thus find the correct value of a . In order to overcome the difficulty, one has to use probabilistic encryption, where each plaintext has many corresponding ciphertexts, depending on some additional random parameter chosen at encryption time. Such a scheme should make it impossible to distinguish encryptions of distinct values, even if these are restricted to range over a set with only two elements. This very strong requirement has been termed *semantic security* ([12]). As a further difference with RSA, the cryptosystem introduced in this paper, has a very natural probabilistic version, with proven semantic security.

The probabilistic homomorphic encryption schemes proposed so far suffer from a serious drawback: they have very poor bandwidth. Typically, they need something like one kilobit to encrypt just a few bits, which is a quite severe expansion rate. This may be acceptable for election schemes but definitely hampers other applications. The main achievement of the present paper is to reach a significant bandwidth, while keeping the other properties, including semantic security.

Before we turn to the more technical developments of our paper, it is in order to compare it with earlier work: it is indeed the case that the question of finding trapdoors for the discrete logarithm problem has been the subject of many papers. At this point, it is fair to mention that the probabilistic cryptosystem that we propose is actually quite close to the most general case of the homomorphic encryption schemes introduced by Benaloh in his Ph-D thesis [4]. Still, both in this thesis and in the related work ([5, 6, 7]), the security and potential applications are only investigated in a setting where the bandwidth remains small. A more recent paper by Park and Won (see [24]) describes a related probabilistic cryptosystem using a trapdoor based on injecting a single power of a small odd integer into $p-1$ or $q-1$ and proves its security with respect to an *ad hoc* statement. Thus, our paper offers the first thorough discussion of the security of a probabilistic homomorphic encryption scheme with significant bandwidth. After the completion of the present work, we have been informed that another homomorphic probabilistic encryption scheme, using moduli n of the form p^2q , where p and q are primes, had been found by Okamoto and Uchiyama (see [22]), achieving an expansion rate similar to ours. Finally, it should be emphasized that the deterministic version of our scheme is not simply a twist that fixes the random string in the probabilistic version: considering its practicality, we believe that, even if it is not intended to be a direct competitor to RSA, it enters the very limited list of efficient public-key cryptosystems.

The paper is organized as follows: in the next two sections, we successively describe the deterministic and the probabilistic version of our scheme, the former with a practical approach, the latter in a more complexity-theoretic spirit. We then discuss applications and end up with a challenge for the research community.

2 The deterministic version

As was just mentioned, our approach to the deterministic scheme is practically oriented: we discuss system set-up

and key-generation, encryption and decryption, with performances in mind. We also carry on a security analysis at the informal level and we derive minimal suggested parameters.

2.1 System set-up and key generation

The scheme that we propose in the present paper can be described as follows: let σ be a squarefree odd B -smooth integer, where B is small integer and let $n = pq$ be an RSA modulus such that σ divides $\phi(n)$ and is prime to $\phi(n)/\sigma$. Typically, we think of B as being a 10 bit integer and we consider n to be at least 768 bits long. Let g be an element whose multiplicative order modulo n is a large multiple of σ . Publish n , g and keep p , q and optionally σ secret. A message m smaller than σ is encrypted by $g^m \bmod n$; decryption is performed using the prime factors of σ as will be seen in the next subsection.

Generation of the modulus appears rather straightforward: pick a family p_i of k small odd distinct primes, with k even. Set $u = \prod_{i=1}^{k/2} p_i$, $v = \prod_{i=k/2+1}^k p_i$ and $\sigma = uv = \prod_{i=1}^k p_i$. Pick two large primes a and b such that both $p = 2au + 1$ and $q = 2bv + 1$ are prime and let $n = pq$.

However, this generation is lengthy especially when the size of the modulus grows: a has to be chosen in the appropriate range and tested for primality as well as $p = 2au + 1$ until both tests succeed simultaneously. This might be a bit time-consuming. Instead, we suggest to generate a , b , u and v first (independently of any primality requirements on p and q) and use a couple of 24-bit "tuning primes" p' and q' (not used in the encryption process) such that $p = 2aup' + 1$ and $q = 2bvq' + 1$ are primes. To avoid interferences with the encryption mechanics, we recommend to make sure that $\gcd(p'q', \sigma) = 1$ and $p' \neq q'$. In practice, such an approach is only 9% slower than equivalent-size RSA key-generation.

To select g , one can choose it at random and check whether or it has order $\phi(n)/4$. The main point is to ensure that g is not a p_i -th power, for each $i \leq k$ by testing that $g^{n/p_i} \not\equiv 1 \bmod n$. The success probability is :

$$\pi = \prod_{i=1}^k \left(1 - \frac{1}{p_i}\right), \text{ whose logarithm is : } \ln(\pi) \approx - \sum_{i=1}^k \frac{1}{p_i}$$

If the p_i s are the first k primes, this in turn can be estimated as $-\ln \ln k$ and results in the quite acceptable overall probability of $\pi \approx 1/\ln k$. Another method consists in choosing, for each index $i \leq k$, a random g_i until it is not a p_i -th power. With overwhelming probability $g = \prod_{i=1}^k g_i^{n/p_i}$ has order $\geq \phi(n)/4$.

2.2 Encryption and Decryption

Encryption consists in a single modular exponentiation: a message m smaller than σ is encrypted by $g^m \bmod n$. Note that it does not require knowledge of σ . A lower bound (preferably a power of two) is enough but it is unclear how important for the security of the scheme is keeping σ secret. However, if one chooses to keep σ secret, necessary precautions (similar to these applied to Rabin's scheme [31] or Shamir's RSA for paranoids [34]) should be enforced for not being used as an oracle¹.

¹For example, an attacker having access to a decryption box can decrypt $g^m \bmod n$ for some $m > \sigma$ and get $m \bmod \sigma$. This discloses (by subtraction) a multiple of σ and σ can then be found by a few re-

Also, there is actually no reason why the p_i s should be prime. Everything goes through, *mutatis mutandis*, as soon as the p_i s are mutually prime. Thus, for example, they can be chosen as prime powers, which is a way to increase the variability of the scheme.

Decryption is based on the chinese remainder theorem. Let $p_i, 1 \leq i \leq k$, be the prime factors of n . The algorithm computes the value m_i of m modulo each p_i and gets the result by chinese remaindering, following an idea which goes back to the Pohlig-Hellman paper [27]. In order to find m_i , given the ciphertext $c = g^m \bmod n$, the algorithm computes $c_i = c^{d(n)/p_i} \bmod n$, which is exactly $g^{m_i d(n)/p_i} \bmod n$. This follows from the following easy computations, where y_i stands for $\frac{d(n)-m_i}{p_i}$:

$$\begin{aligned} c_i &= c^{d(n)/p_i} = g^{m d(n)/p_i} = g^{(m_i + y_i p_i) d(n)/p_i} \\ &= g^{m_i d(n)/p_i} g^{y_i d(n)} = g^{m_i d(n)/p_i} \bmod n \end{aligned}$$

By comparing this result with all possible powers $g^{j d(n)/p_i}$, it finds out the correct value of m_i . In other words, one loops for $j = 0$ to $p_i - 1$ until $c_i = g^{j d(n)/p_i} \bmod n$.

The cleartext m can therefore be computed by the following procedure:

```

for i = 1 to k
{
  let  $c_i = c^{d(n)/p_i} \bmod n$ 
  for j = 0 to  $p_i - 1$ 
  {if  $c_i = g^{j d(n)/p_i} \bmod n$  let  $m_i = j$ }
}
x = ChineseRemainder({ $m_i$ }, { $p_i$ })

```

The basic operation used by this (non-optimized) algorithm is a modular exponentiation of complexity $\log^3(n)$, repeated less than:

$$k p_k < \log(n) p_k \leq \log(n) k \log(k) < \log^2(n) \log \log(n)$$

times. Decryption therefore takes $\log^5(n) \log \log(n)$ bit operations.

This is clearly worse than the $\log^3(n)$ complexity of RSA but encryption can be optimized if a table stores all possible values of $t[i, j] = g^{j d(n)/p_i} \bmod n$, for $1 \leq i \leq k$ and $1 \leq j \leq p_i$: the value m_i of the cleartext m modulo p_i is found by table look-up, once $c^{d(n)/p_i} \bmod n$ has been computed. It is not really necessary to store all $g^{j d(n)/p_i}$. Any hash function that distinguishes $g^{j d(n)/p_i}$ from $g^{j' d(n)/p_i}$, for $j \neq j'$ will do and, in practical terms, a few bytes will be enough, for example approximately $2|p_i|$ bits from each $t[i, j]$. It is even possible to use hash functions that do not discriminate values of $g^{j d(n)/p_i}$: the proper one is spotted by considering, by table look-up

peated trials and gcda. To prevent such an action, the decryption box cannot only re-encrypt and check against the ciphertext received, as this allows a search by dichotomy. It should first check that the cleartext is in the appropriate range, e.g. $< 2^l$ with $2^l < m$, re-encrypt it and then check that it matches up with the original ciphertext before letting anything out.

hashes of $g^{j d(n)/p_i}$, for $\ell = 1, 2, \dots$ until there is no ambiguity. This can be very efficiently implemented by storing hash values in increasing order w.r.t. ℓ and one single bit might be enough.

2.3 A toy example

• key generation for $k = 6$

$$p = 21211 = 2 \times 101 \times 3 \times 5 \times 7 + 1,$$

$$q = 928643 = 2 \times 191 \times 11 \times 13 \times 17 + 1,$$

$n = 21211 \times 928643 = 19697446673$ and $g = 131$ yield the table:

| | $i=1$ | $i=2$ | $i=3$ | $i=4$ | $i=5$ | $i=6$ |
|--------|-------|-------|-------|-------|-------|-------|
| $j=0$ | 0001 | 0001 | 0001 | 0001 | 0001 | 0001 |
| $j=1$ | 1966 | 6544 | 1967 | 6273 | 6043 | 0372 |
| $j=2$ | 9560 | 3339 | 4968 | 7876 | 4792 | 7767 |
| $j=3$ | | 9400 | 1765 | 8720 | 0262 | 3397 |
| $j=4$ | | 5479 | 6701 | 7994 | 0136 | 0702 |
| $j=5$ | | | 6488 | 8651 | 6291 | 4586 |
| $j=6$ | | | 2782 | 4691 | 0677 | 8135 |
| $j=7$ | | | | 9489 | 1890 | 3902 |
| $j=8$ | | | | 8537 | 6878 | 5930 |
| $j=9$ | | | | 2312 | 2571 | 6399 |
| $j=10$ | | | | 7707 | 7180 | 6592 |
| $j=11$ | | | | | 8291 | 9771 |
| $j=12$ | | | | | 0678 | 0609 |
| $j=13$ | | | | | | 7337 |
| $j=14$ | | | | | | 6892 |
| $j=15$ | | | | | | 3370 |
| $j=16$ | | | | | | 3489 |

where entry $\{i, j\}$ contains $g^{j d(n)/p_i} \bmod n \bmod 10000$.
• encryption of $m = 202$

$$c = g^m \bmod n = 131^{202} \bmod 19697446673 = 519690214$$

• decryption
by exponentiation, we retrieve:

$$\begin{aligned} c^{d(n)/p_1} \bmod n \bmod 10000 &= 1966 \\ c^{d(n)/p_2} \bmod n \bmod 10000 &= 3339 \\ c^{d(n)/p_3} \bmod n \bmod 10000 &= 2782 \\ c^{d(n)/p_4} \bmod n \bmod 10000 &= 7994 \\ c^{d(n)/p_5} \bmod n \bmod 10000 &= 1890 \\ c^{d(n)/p_6} \bmod n \bmod 10000 &= 3370 \end{aligned}$$

wherefrom, by table lookup:

$$\begin{aligned} m \bmod 3 &= \text{table}(1966) = 1 \\ m \bmod 5 &= \text{table}(3339) = 2 \\ m \bmod 7 &= \text{table}(2782) = 6 \\ m \bmod 11 &= \text{table}(7994) = 4 \\ m \bmod 13 &= \text{table}(1890) = 7 \\ m \bmod 17 &= \text{table}(3370) = 15 \end{aligned}$$

and by Chinese remaindering: $m = 202$.

2.4 Suggested parameters and security analysis

We suggest to take $\sigma > 2^{160}$ and we consider $|n| = 768$ bits as a minimum size for the modulus.

If the factorization of n is found, then a and b become known as well as $\phi(n)$. The scheme is therefore broken. However, the scheme does not appear to be provably equivalent to factoring. Rather, it is related to the question of having oracles that decide whether or not a random number x is a p_i -th power modulo n , for $i = 1, \dots, k$. This is known as the higher residuosity problem and is currently considered unfeasible. Formal equivalence of this problem and the probabilistic version of our encryption scheme will be proved in the next session. Considering the basic deterministic version, we have no formal proof but we haven't found any plausible line of attack either. Also, the efficient factoring methods such as the quadratic sieve (QS) or the number field sieve (NFS) do not appear to take any advantage from the side information that u (resp. v) divides $p-1$ (resp. $q-1$). The same is true of simpler methods like Pollard's $p-1$ since we have ensured that neither $p-1$ nor $q-1$ is smooth. Finally, elliptic curve weaponry [18] will not pull-out factors of n in the range considered. Note that the requested size of n (768 bits or more) makes factoring n a very hard task anyway.

We now turn the size of σ . In order to avoid the computation of discrete logarithms by the baby step-giant step method, we have to make σ large enough. As already stated, 2^{160} is a minimum. This can be achieved for example by making σ a permutation of the first 30 odd primes, which yields $\sigma \approx 2^{160.45}$. Alternatively, one can choose a sequence of 16 primes with 10 bits. Since there are 75 such primes, this leads to a ≈ 53 -bit entropy. Adding prime powers, as stated above, will further increase these figures.

There is a further difficulty, when σ is known. Note that

$$4ab = \frac{\phi(n)}{\prod_{i=1}^k p_i} = \frac{n - p - q + 1}{\sigma}$$

hence $4ab$ differs from $\frac{n}{\sigma}$ only by $\epsilon = -\frac{p+q-1}{\sigma}$. The numerator is of size $|n|/2$, hence, if it does not exceed the denominator by a fairly large number of bits, the value of ab is basically known and decryption can be performed.

When the exact splitting of the factors of σ into u and v are known as well, the previous analysis can be pushed further. Reducing the relation $n = (2au+1)(2bv+1)$ modulo u , we find that $n = 2bv + 1 \pmod{u}$ and we can calculate $d = b \pmod{u}$. Similarly, we learn $c = a \pmod{v}$. We let $a = rv + c$ and $b = su + d$, with r, s unknown and, using the fact that $\sigma = uv$, we obtain:

$$n = (2rvu + 2cu + 1)(2su + 2dv + 1) =$$

$$4rs\sigma^2 + 2\sigma[r(2dv + 1) + s(2cu + 1)] + (2cu + 1)(2dv + 1)$$

which is of the form

$$n = 4rs\sigma^2 + 2\sigma(\alpha r + \beta s) + \gamma$$

with known α, β and γ . Reducing modulo σ^2 , this provides the value δ of $\alpha r + \beta s \pmod{\sigma}$. At this point, our analysis becomes quite technical and the reader may skip the following and jump to the conclusion that $n \gg \sigma^4$.

For the interested reader, we note that the pair (r, s) lies in the two-dimensional lattice L defined by

$$L = \{(x, y) | \alpha x + \beta y = \delta \pmod{\sigma}\}$$

This lattice has determinant σ . Also, it is easily seen that α and β are bounded by 2σ and γ by $4\sigma^2$. From this we get

$$rs \leq \frac{n}{4\sigma^2} \leq rs + r + s + 1 = (r+1)(s+1)$$

Thus, the pair (r, s) is very close to the boundary of the curve C with equation $xy = \frac{n}{4\sigma^2}$. More precisely, the distance between the pair (r, s) and the curve does not exceed $\sqrt{2}$. This defines a geometric area A that includes (r, s) . Now, key generation usually induces constraints that limit the possible range of the parameters. For this reason, it is appropriate to replace C by the line $x + y = \frac{\sqrt{n}}{\sigma}$ in order to estimate the size of A . This leads to an approximation which is $O(\frac{\sqrt{n}}{\sigma})$. The number of lattice points from L in this area is, in turn, measured by the ratio between the size of A and the determinant, which is $\frac{\sqrt{n}}{\sigma^2}$. It is safe to ensure that this set is beyond exhaustive search, which we express by $n \gg \sigma^4$.

Note that the ratio $|n|/|\sigma|$ is the expansion rate of the encryption, where $|n|$ denotes, as usual, the size of n in bits. It is of course desirable to make this rate as low as possible. On the other hand, as a consequence of the above remarks, we see that $\frac{|n|}{|\sigma|} - |\sigma|$ should be large. Asymptotically, this is achieved as soon as we fix an expansion rate which is > 4 . For real-size parameters, we suggest to respect the heuristic bound $\frac{|n|}{|\sigma|} - \sigma \geq 128$, which is consistent with our minimal parameters. Larger parameters allow a slightly better expansion rate.

2.5 Performances

Despite its expansion rate, the new cryptosystem is quite efficient: encryption requires the elevation of a constant 768-bit number to a 160-bit power. Several batch ([21, 23]) and pre-processing ([2]) techniques can speed-up such computations, which might be a small advantage over RSA.

Decryption is slightly more awkward since k exponentiations are needed. But this number can be reduced in a few ways:

Firstly, while computing $c^{\phi(n)/p_i} \pmod{n}$ for each i , it is possible to first store $c' = c^{4ab} \pmod{n}$ and raise c' to the successive powers σ/p_i so that (besides the first one), the remaining exponentiations involve 160-bit powers. One can further, in the square-and-multiply algorithm, share the "square" part of the various exponentiations. A careful bookkeeping of the number of modular multiplications obtained by setting $|n| = 768$ and choosing sixteen 10-bit primes p_i , shows that the total number of modular multiplications decreases to 2352: 912 for the computation of c' and 1440 for the rest. Actually, the "multiply" part can be somehow amortized as well: we refer to [21] for a proper description of such an optimized exponentiation strategy. The resulting computing load is less than what is needed for a couple of RSA decryptions with a similar modulus.

Unfortunately, there is a drawback in reducing the value of k : in the 30-prime variant it is necessary to store 1718 different $\{i, j\}$ hash values. Hashing on two bytes seems enough and results in an overall memory requirement of four kilobytes. In the 16-prime variant, hash values of 3 bytes are necessary and the table size becomes ≈ 100 kilobytes. As observed at the end of section 2.2, the hash table can be drastically reduced at the cost of a minute computation overhead.

Another speed-up can be obtained by separately performing decryption modulo p and q so as to take advantage

of smaller operand sizes. This alone, divides the decryption workload by four.

Finally, decryption is inherently parallel and naturally adapted to array processors since each m_i can be computed independently of all the others.

2.6 Implementation

The new scheme (768-bit n , $k = 30$) was actually implemented on a 68HC05-based ST16CF54 smart-card (4,096 EEPROM bytes, 16,384 ROM bytes and 352 RAM bytes). The public key is only 96-byte long and as in most smart-card implementations, n 's storage is avoided by a command that re-computes the modulus from its factors upon request (re-computation and transmission take 10 ms). For further space optimization g 's first 91 bytes are the byte-reversed binary complement of n 's last 91 bytes. Decryption (a 4,119-byte routine) takes 3,912 ms. Benchmarks were done with a 5 MHz oscillator and ISO 7816-3 T=0 transmission at 115,200 bauds.

3 The probabilistic version

3.1 The setting

We now turn to the probabilistic version of the scheme. As already explained, we adopt a more complexity-oriented approach and, for example, we view B as bounded by a polynomial in $\log n$. The probabilistic version replaces the ciphertext $g^m \bmod n$ by $c = x^m g^m \bmod n$, where x is chosen at random among positive integers $< n$. Decryption remains identical. This is due to the fact that the effect of multiplying by x^p is cancelled by raising the ciphertext to the various powers $\frac{\phi(n)}{p}$, as performed by the decryption algorithm. Note that this version requires σ to be public.

The resulting scheme is *homomorphic*, which means that $E(m + m' \bmod \sigma) = E(m)E(m') \bmod n$. Probabilistic homomorphic encryption has received a lot of applications, both practically and theoretically oriented. To name a few, we quote the early work of Benaloh on election schemes ([4]) and the area of zero-knowledge proofs for NP (see [13, 3]). Known such schemes are the Quadratic Residuosity schemes of Goldwasser and Micali ([12]) which encrypts only one bit and its extensions to higher residues modulo a single prime (see [4]), which encrypts a few bits. As already explained in section 1, these schemes suffer from a serious drawback: a complexity theoretic analysis has to view the cleartext as logarithmic in the size of ciphertext. In other words, the expansion rate, i.e. the ratio between the length of the ciphertext and the length of the cleartext is huge. In our proposal, this ratio is exactly $\frac{|n|}{|\sigma|}$. Note that that our assumption that σ is B -smooth, for some small B , does not preclude a linear ratio. The maximum size of σ is $\sum_{p \leq B} \log p$, where p ranges over primes and it is known that $\theta(B) = \sum_{p \leq B} \ln p \simeq B$. Thus, even if B is logarithmic in n , there are enough primes to make $|\sigma|$ a linear proportion of $|n|$. This is a definite improvement over previous homomorphic schemes. Note however that, following the comments in section 2.4, it is safe to take $\frac{|\sigma|}{|n|} < 1/4$.

3.2 A complexity theoretic approach

We already observed that the security of our proposal is related to the question of distinguishing higher residues modulo n , that is integers of the form $x^p \bmod n$, when p is a

prime divisor of $\phi(n)$. In the rest of this section, we want to clarify this relationship in the asymptotic setting of complexity theory. In view of the remarks just made, we find it convenient to assume that the ratio $\frac{|\sigma|}{|n|}$ has a fixed value $\alpha < 1/4$. We also fix a polynomial B in $\log n$. The parameters which are of interest to us are pairs (n, σ) such that σ is squarefree, odd and B -smooth, n is a product of two primes p, q , σ is a divisor of $\phi(n)$ prime to $\phi(n)/\sigma$ and $\frac{|\sigma|}{|n|} = \alpha$. We call any integer n that appears as first coordinate of such a pair (B, α) -dense. Distinguishing higher residues is usually considered difficult (see [4]). We conjecture that this remains true when n varies over (B, α) -dense integers. Towards a more precise statement, let $R_p(y, n)$ be one if y is a p -th residue modulo n and zero otherwise. Define a higher residue oracle to be a probabilistic polynomial time algorithm A which takes as input a triple (n, y, p) and returns a bit $A(n, y, p)$ such that the following holds: There exists a polynomial Q in $|n|$ such that, for infinitely many values of $|n|$, one can find a prime $p(|n|) < B$, with:

$$\Pr\{A(n, y, p) = R_p(y, n)\} \geq 1 - \frac{1}{p} + \frac{1}{Q}$$

where the probability is taken over the random tosses of A and its inputs, conditionally to the event that n is (B, α) -dense and p is a divisor of $\phi(n)$.

Our Intractability Hypothesis is that there is no higher residue oracle. The constant $1 - \frac{1}{p}$ comes from the obvious strategy for approximating R_p which consists in constantly outputting zero. This strategy is successful for a proportion $1 - \frac{1}{p}$ of the inputs.

3.3 A security proof

The security of probabilistic encryption scheme has been investigated in [12]. In this paper, the authors introduced the notion of *semantic security*: given two messages m_0 and m_1 , a message distinguisher is a probabilistic polynomial time algorithm D , which distinguishes encryptions of m_0 from encryptions of m_1 . More, accurately, it outputs a bit $D(n, \sigma, g, y)$ in such a way that, setting

$$\theta_i = \Pr\{D(n, \sigma, g, y) = 1 | y \in E(m_i)\}$$

where $E(m_i)$ is the set of encryptions of m_i , the following holds:

There exists a polynomial Q in $|n|$ such that, for infinitely many values of $|n|$, $|\theta_0 - \theta_1| \geq \frac{1}{Q}$.

Semantic security is the assertion that there is no pair of polynomial time algorithms F, D such that F produces two messages for which D is a message distinguisher.

Theorem 1 Assume that no higher residue oracle exists. Then, the probabilistic version of the encryption scheme has semantic security.

The proof of this result uses the *hybrid technique* for which we refer to [11]. It is technical in character and we have chosen to only include a sketch it in an appendix to the present paper.

4 Applications and variants

Even if we do not expect large scale replacement of RSA by our scheme, we feel that the latter is worth some academic interest. Especially, we believe that it opens up new applications. We have not yet fully investigated those potential applications but we give some suggestions below.

4.1 Traceability

Our proposal could offer some help in the management of key escrowing services. Consider the variant of the Diffie-Hellman key exchange protocol, where a composite modulus n is used. Such a variant has been studied by various researchers including Mc Curley in [20], where it is shown that some specific choices lead to a scheme that is at least as difficult as factoring. Assume further that the modulus n and the base for exponentiations g are chosen as described in section 1. It has been proposed (see e.g. [14]) that g and n could be defined by some kind of TTP (Trusted Third Party). Now, the user's public key y and his secret key x are related by $y = g^x \bmod n$. It is conceivable to leave the choice of x to the user with the provision that $x \bmod \sigma = ID$, where ID is the identity of the user. This can be checked by the TTP upon registration of the key. Thus, we have reached a situation where the identity is embedded in the public key through a trapdoor, although the actual key is not. One should not however overestimate the resulting functionality. It could be useful in scenarios where traceability is made possible via escrowing but where confidentiality cannot be broken even with the help of the escrowing services. Alternatively, it might be used to split traceability and secret key recovery between key escrows. Note that the above proposal requests that σ is made public: as already observed, this does not seem to endanger the scheme.

4.2 Variants of the scheme

As is often the case, one can design numerous variants of the basic scheme. We will mention two because of their potential applications.

Use of moduli with three prime factors As for RSA, it is possible to embed three prime factors p, q, r in the modulus in place of two. The construction is straightforward: the small odd primes p_i are split into three groups thus yielding, by multiplication, three integers u, v, w . The three primes are then sought among integers of the form $2au + 1$ (resp. $2bv + 1$, resp. $2cw + 1$). It seems possible to keep the minimum size of n to 768 bits, which allows a, b, c to be around 200 bits. Following an idea of Maurer and Yacobi ([19]), we can then have a complete trapdoor for the discrete logarithm with base g : once the σ part has been computed, there remains to compute the logarithm modulo a, b and c , which is not immediate but well within the reach of current technology, since these numbers are 200 bit integers. Again, the variant could prove useful in key escrowing scenarios of, say, Diffie-Hellman keys, where it might be desirable to have a lengthy recovery of the secret key for consumer's protection.

Multiplicative encryption In this variant, σ is made public and encryption applies to messages of length k , $m = \sum_{i=1}^k m_i 2^{i-1}$. In order to encrypt m , one computes $e = \prod_{i=1}^k p_i^{m_i}$ and apply probabilistic encryption to e . Of course, the bandwidth of this variant is very low: using a 768 bit modulus n and choosing the first 30 odd primes for p_i s, we obtain a 30 bit input and a 768 bit output. Allowing a larger input has drastic consequences in terms of the size of n . The value of σ is close to $2^{k\phi}$ when the first k primes are used with $k = 80$ but reaches $2^{928.4}$ for $k = 128$ and 2^{1309} for $k = 160$. Using the heuristic bound mentioned in section 2.4, we get for the length of n something beyond 5000 bits if k is 160. This goes down to 2400 bits when $k = 80$.

As a result, the variant just described is not really practical and there is little chance that it can ever be adopted as an actual encryption scheme. On the other hand, the ciphertext $c(m)$ can be used in an encryption scheme à la El Gamal. The modulus is not prime since it is an RSA modulus, but it makes no difference on the user's size. From $h = c(m)$, he can manufacture a public key y with a corresponding matching secret key x of his choice $y = h^x \bmod n$. The resulting cryptosystem allows ciphertext traceability in the sense of Desmedt (see [9]). Our proposal enables to trace ciphertexts by a technique similar to the one used by Desmedt, but decreases the size of the modulus from something like 10000 bits to 2500 bits. The tracing algorithm goes as follows: extract from an El Gamal encryption the part $v = h^r \bmod n$ and apply the decryption algorithm, treating v as a ciphertext. The decryption algorithm will basically find the original message m , which provides the identity of the user and from which h was built. Several errors may occur due to the fact that r might have some of the p_i s as divisors: the corresponding decrypted values of m_i will be set to 1, regardless of their original values. The correct value can be found if a sample of ciphertexts are available or, alternatively, if an error-correction capacity has been added to m . Such an error-correction mechanism is highly advisable anyway in view of the attacks against software key escrow reported in [15].

Note that, one can further reduce the size of the exponent. This is because 40 bits may be considered enough for tracing purposes. The value of σ goes down to approximately 2^{243} and 1088 bits becomes an acceptable minimum length for the modulus.

5 Challenge

It is a tradition in the cryptographic community to offer cash rewards for successful cryptanalysis. More than a simple motivation means, such rewards also express the designers' confidence in their own schemes. As an incentive to the analysis of the new scheme, we therefore offer \$ $|n|$ to whoever will decrypt :

```
c = 13370fe62d81fde356d1842fd7e5fc1ae5b9b449
b4d00866597e61af4fb0d939283b04d3bb73f91f
0d9d61eb0014690e567ab89aa8df4a9164cd4c6e
6df80806c7cdeda5cfd97bf7c42cc702512a49
dd196c8746c0e2ef36ca2ae21d4a36a1e

g = 0b9cf6a789959ed4f36b701a5065154f7f4f1517
6d731b4897875d26a9e24415e111479050894ba7
c532ada1903c63a84ef7edc29c208a8dd43fb5f7
d43727b730f20d8e12c17cd5cf9ab4358147cb62
a9fb8878bf15204e444ba6ade61327431e

n = 1459b9617b8a9df6bd54341307f1256dafa241bd
65b96ed14078e80dc6116001b83c5f88c7bbcb0b
db237daac2e76df5b415d089baa0fd078516e60e
2cdda7c26b858777604c5fbd19f0711bc75ce00a
5c37e2790b0d9d0ff9625c5ab9c7511d1e
```

where $k = 30$ (p_i is the i -th odd prime) and the message is ASCII-encoded. The challenger should be the first to decrypt at least 50% of c and publish the cryptanalysis method but the authors are ready to carefully evaluate ad valorem any feedback they get.

Acknowledgements

The paper grew out of a previous version which did not include the probabilistic case of our scheme. We wish to thank Julien Stern for suggesting us this alternative mode of encryption. We also want to thank J. Benaloh for help in clarifying our respective contributions in the definition of the probabilistic case. Finally, we are grateful to Adi Shamir, for helpful comments including the improved decryption algorithm mentioned in section 2.2 and also to one of the anonymous referees for pointing out the clever trick that yields the improved security analysis included at the end of section 2.4.

References

- [1] R. Anderson, *Robustness principles for public-key protocols*, Advances in Cryptology Crypto'95, Santa Barbara, Lectures Notes in Computer Science 963, pp. 236-247, Springer-Verlag, 1995.
- [2] E. Brickell, D. Gordon, K. McCurley and D. Wilson, *Fast Exponentiation with Precomputation*, Advances in Cryptology Eurocrypt'92, Balatonfured, Lectures Notes in Computer Science 658, pp. 200-207, Springer-Verlag, 1993.
- [3] G. Brassard, D. Chaum and C. Crépeau, *Minimum Disclosure Proofs of Knowledge*, JCSS, Vol. 37(2), Oct. 1988, pp. 156-189.
- [4] J. D. Cohen Benaloh, *Verifiable Secret-Ballot Elections*, Ph-D thesis, Yale University, 1988.
- [5] J. D. Cohen and M. J. Fischer, (1985), *A robust and verifiable cryptographically secure election scheme*, Proc. of 26th Symp. on Foundation of Computer Science, 1985, 372-382.
- [6] J. D. Cohen Benaloh, *Cryptographic Capsules: A Disjunctive Primitive for Interactive Protocols*, Advances in Cryptology Crypto'86, Santa Barbara, Lectures Notes in Computer Science, pp. 213-222, Springer-Verlag, 1986.
- [7] J. D. Cohen Benaloh and M. Yung, *Distributing the Power of a Government to Enhance the Privacy of Voters*, Proc. of 5th Symp. on Principles of Distributed Computing, 1986, 52-62.
- [8] D. Denning (Robling), *Cryptography and data security*, Addison-Wesley Publishing Company, pp. 148, 1983.
- [9] Y. Desmedt, *Securing traceability of ciphertexts - Towards a secure software key escrow system*, Advances in Cryptology Eurocrypt'95, Saint-Malo, Lectures Notes in Computer Science 921, pp. 417-457, Springer-Verlag, 1995.
- [10] W. Diffie and M. Hellman, *New directions in cryptography*, IEEE Transactions on Information Theory, vol. IT-22-6, pp. 644-654, 1976.
- [11] O. Goldreich, *Foundations of cryptography (Fragments of a book)*, Weizmann Institut of Science, 1995.
- [12] S. Goldwasser and S. Micali, *Probabilistic Encryption*, JCSS, 28(2), April 1984, pp. 270-299.
- [13] O. Goldreich, S. Micali and A. Wigderson, *Proofs that Yield Nothing but their Validity and a Methodology of Cryptographic Protocol Design*, Proc. of 27th Symp. on Foundation of Computer Science, 1986, pp.174-187.
- [14] N. Jefferies, C. Mitchell and M. Walker, *A proposed architecture for trusted third party services*, Cryptography Policy and Algorithms, Queensland, Lecture Notes in Computer Science 1029, pp. 98-114, Springer-Verlag, 1996.
- [15] L. Knudsen and T. Pedersen, *On the difficulty of software key escrow*, Advances in Cryptology Eurocrypt'96, Saragossa, Lectures Notes in Computer Science 1070, pp. 237-244, Springer-Verlag, 1996.
- [16] P. Kocher, *Timing attacks in implementations of Diffie-Hellman, RSA, DSS and other systems*, Advances in Cryptology Crypto'96, Santa Barbara, Lectures Notes in Computer Science, pp. 104-113, Springer-Verlag, 1996.
- [17] Kaoru Kurosawa, Yutaka Katayama, Wakaha Ogata and Shigeo Tsujii, *General public key residue cryptosystems and mental poker protocols*, Advances in Cryptology Eurocrypt'90, Aarhus, Lectures Notes in Computer Science 473, pp. 374-388, Springer-Verlag, 1996.
- [18] H. Lenstra Jr., *Factoring integers with elliptic curves*, Annals of Mathematics, 126, pp. 649-673, 1991.
- [19] U. Maurer and Y. Yacobi, *Non-interactive public key cryptography*, Advances in Cryptology Eurocrypt'91, Brighton, Lectures Notes in Computer Science 547, pp. 498-507, Springer-Verlag, 1991.
- [20] K. McCurley, *A key distribution system equivalent to factoring*, Journal of Cryptology, vol. 1, pp. 85-105, 1988.
- [21] D. M'Raihi and D. Naccache, *Batch exponentiation - A fast DLP-based signature generation strategy*, Proceedings of the third ACM conference on Computer and Communications Security, New Delhi, pp. 58-61, 1996.
- [22] T. Okamoto and S. Uchiyama, *A new public-key cryptosystem as secure as factoring*, Advances in Cryptology Eurocrypt'98, Helsinki, Lectures Notes in Computer Science, pp. to appear, Springer-Verlag, 1998.
- [23] D. Naccache and J. Stern, *A new public-key cryptosystem*, Advances in Cryptology Eurocrypt'97, Constance, Lectures Notes in Computer Science 1233, pp. 27-36, Springer-Verlag, 1997.
- [24] Sung-Jun Park and Dong-Ho Won, *A generalization of public key residue cryptosystem*, In Proc. of 1993 KOREA-JAPAN joint workshop on information security and cryptology, 202-206.
- [25] B. Pfitzmann and M. Schunter, *Asymmetric fingerprinting*, Advances in Cryptology Eurocrypt'96, Saragossa, Lectures Notes in Computer Science 1070, pp. 84-95, Springer-Verlag, 1996.

- [26] D. Pointcheval, *A new identification scheme based on the perceptrons problem*, Advances in Cryptology Eurocrypt'94, Perugia, Lecture Notes in Computer Science 950, pp. 318-328, Springer-Verlag, 1995.
- [27] S. C. Pohlig and M. E. Hellman, *An improved algorithm for computing logarithms over $GF(p)$ and its cryptographic significance* IEEE Transactions on Information Theory, vol. IT-24-1, pp. 105-110, 1978.
- [28] J. Pollard, *Theorems on factorization and primality testing*, Proceedings of the Cambridge Philosophical Society, vol. 76, pp. 521-528, 1974.
- [29] J. Pollard, *Factoring with cubic integers*, A. Lenstra and H. Lenstra Jr., The development of the number field sieve, vol. 1554, LNM, 4-10, Springer-Verlag, 1993.
- [30] C. Pomerance, *Analysis and comparison of some integer factoring algorithms*, printed in H. Lenstra Jr. and R. Tijdeman, Computational Methods in Number Theory I, Mathematisch Centrum Tract 154, Amsterdam, pp. 89-139, 1982.
- [31] M. Rabin, *Digitalized signatures and public-key functions as intractable as factorization*, MIT/LCS/TR-212, MIT Laboratory for Computer Science, 1979.
- [32] R. Rivest, A. Shamir and L. Adleman, *A method for obtaining digital signatures and public-key cryptosystems*, Communications of the ACM, vol. 21-2, pp. 120-126, 1978.
- [33] A. Shamir, *An efficient identification scheme based on permuted kernels*, Advances in Cryptology Crypto'89, Santa Barbara, Lecture Notes in Computer Science 435, pp. 606-609, Springer-Verlag, 1990.
- [34] A. Shamir, *RSA for paranoids*, CryptoBytes, vol. 1-3, pp. 1-4, 1995.
- [35] J. Stern, *A new identification scheme based on syndrome decoding*, Advances in Cryptology Crypto'93, Santa Barbara, Lecture Notes in Computer Science 773, pp. 13-21, Springer-Verlag, 1994.
- [36] J. Stern, *Designing identification schemes with keys of short size*, Advances in Cryptology Crypto'94, Santa Barbara, Lecture Notes in Computer Science 839, pp. 164-173, Springer-Verlag, 1995.

Appendix: Sketch of the Security Proof.

We show that any message distinguisher can be turned into an algorithm that recognizes higher residues. We let D be a distinguisher for two messages m_0 and m_1 and start from the fact that, keeping the above notations, θ_0 and θ_1 are significantly distinct. We next use the *hybrid technique* for which we refer to [11], pp.91-93. Hybrids consist of a sequence of random variables Y_i , $0 \leq i \leq k$, such that

1. Extreme hybrids collide with $E(m_0)$ and $E(m_1)$ respectively.
2. Random values of each hybrid can be produced by a probabilistic polynomial time algorithm.

3. There are only polynomially many hybrids.

In such a situation, [11] shows that D distinguishes two neighbouring hybrids. Our hybrids are formed by considering a message μ_i , such that

$$\mu_i = m_0 \bmod p_j \text{ for } j > i \text{ and}$$

$$\mu_i = m_1 \bmod p_j \text{ for } j \leq i$$

and letting Y_i to be uniformly distributed over the set $E(\mu_i)$ of encryptions of μ_i . It is easily seen that conditions 1, 2 and 3 are satisfied. Thus, for some index i , D significantly distinguishes Y_i and Y_{i-1} . Set $\mu = \mu_i$, $p = p_i$ and let μ^ℓ , $1 \leq \ell \leq p$, be the unique message such that

$$\mu^\ell = \mu \bmod p_\ell \text{ for } \ell \neq i \text{ and } \mu^\ell = j \bmod p$$

We note that, both m_i and m_{i-1} appear among the μ^ℓ 's and we show that D cannot distinguish encryptions of any two of the μ^ℓ 's. This will yield the desired contradiction.

Let

$$\pi_j = \Pr\{D(n, \sigma, g, y) = 1 | y \in E(\mu_j)\}$$

and assume that some π_i significantly exceeds the other ones. In other words, $\pi_i \geq \sup_{j \neq i} \pi_j + \frac{1}{Q}$ for some polynomial Q and infinitely many values of $|n|$. We show how to predict p -th residuosity: given x , we run D over a large sample N of inputs (n, σ, y) where $y = x^{\sigma} z^{\ell/p} g^{\mu^\ell}$, with $x > n$ and $\ell \leq p$ chosen at random, and we average the outputs. Now, if x is a p -th residue, then y simply varies over $E(\mu_i)$, whereas, if x is not a p -th residue, y randomly varies over the union of all $E(\mu_j)$'s. Thus, in the first case, the average is close to π_i , whereas, in the second case, it is

approximately $\frac{\sum_{j \neq i} \pi_j}{p}$. It is easily seen that the difference is bounded from below by $\frac{p-1}{p} \frac{1}{Q}$. Using the law of large numbers, this is enough to make the proper decision on the p -th residuosity, with probability as close to 1 as we wish, by using only polynomially large samples. This finishes the proof.

Remarks.

1. Turning the previous sketch into a complete proof involves a technical but rather long write-up: especially, a precise version of the law of large numbers has to be made explicit, e.g. by using the Chebishev inequality. Also, the values of π_i and $\frac{\sum_{j \neq i} \pi_j}{p}$ are not known a priori and should be approximated as well using the law of large numbers. We urge the interested reader to consult [11] for similar proofs.
2. The higher residuosity oracle that was built in the proof for the sake of contradiction uses inputs σ and g on top of n , y and p . Actually, one can check that everything goes through, *mutatis mutandis*, if σ is replaced by $\bar{\sigma} = \prod_{p < B} p$. Thus σ is not really needed. As for g , as seen in section 2.1, it can be chosen at random: a proper choice will be spotted by sampling the corresponding oracle and checking its correctness.

Twin Signatures: an Alternative to the Hash-and-Sign Paradigm

✓ David Naccache
Gemplus Card International
34, rue Guynemer
92447 Issy-les-Moulineaux, France
david.naccache@gemplus.com

David Pointcheval Jacques Stern
École Normale Supérieure
45, rue d'Ulm
75230 Paris cedex 05, France
(david.pointcheval,jacques.stern)@ens.fr

ABSTRACT

This paper introduces a simple alternative to the hash-and-sign paradigm, from the security point of view but for signing short messages, called *twinning*. A twin signature is obtained by signing twice a short message by a signature scheme. Analysis of the concept in different settings yields the following results:

- We prove that no generic algorithm can efficiently forge a twin DSA signature. Although generic algorithms offer a less stringent form of security than computational reductions in the standard model, such successful proofs still produce positive evidence in favor of the correctness of the new paradigm.
- We prove in standard model an equivalence between the hardness of producing existential forgeries (even under adaptively chosen message attacks) of a twin version of a signature scheme proposed by Gennaro, Halevi and Rabin and the Flexible RSA Problem.

We consequently regard twinning as an interesting alternative to hash functions for eradicating existential forgery in signature schemes.

Keywords

Digital Signatures, Provable Security, Discrete Logarithm, Generic Model, Flexible RSA Problem, Standard Model.

1. INTRODUCTION

The well-known *hash and sign* paradigm has two distinct goals: increasing *performance* by reducing the size of the signed message and improving *security* by preventing existential forgeries. As a corollary, hashing remains mandatory even for short messages.

From the conceptual standpoint, the use of hash functions comes at the cost of extra assumptions such as the conjecture that for all practical purposes, concrete functions can

be identified with ideal black boxes [3] or that under certain circumstances (black box groups [15, 21]) a new group element must necessarily come from the addition of two *already known* elements. In some settings [11] both models are even used simultaneously.

This paper investigates a simple substitute to hashing that we call *twinning*. A twin signature is obtained by signing twice the same (short) raw message by a probabilistic signature scheme, or two probabilistically related messages.

We believe that this simple paradigm is powerful enough to eradicate existential forgery in a variety of contexts. To support this claim, we show that no generic algorithm can efficiently forge a twin DSA signature and prove that for a twin variant of a signature scheme proposed by Gennaro, Halevi and Rabin [8] (hereafter GHR) existential forgery, even under an adaptively chosen-message attack, is equivalent to the Flexible RSA Problem [5] in the standard model.

2. DIGITAL SIGNATURE SCHEMES

Let us begin with a quick review of definitions and security notions for digital signatures. Digital signature schemes are the electronic version of handwritten signatures for digital documents: a user's signature on a message m is a string which depends on m , on public and secret data specific to the user and—possibly—on randomly chosen data, in such a way that anyone can check the validity of the signature by using public data only. The user's public data are called the *public key*, whereas his secret data are called the *secret key*. The intuitive security notion would be the impossibility to forge user's signatures without the knowledge of his secret key. In this section, we give a more precise definition of signature schemes and of the possible attacks against them (most of those definitions are based on [9]).

2.1 Definitions

A signature scheme is defined by the three following algorithms:

- The *key generation algorithm* G . On input 1^k , where k is the security parameter, the algorithm G produces a pair (k_p, k_s) of matching public and secret keys. Algorithm G is probabilistic.
- The *signing algorithm* Σ . Given a message m and a pair of matching public and secret keys (k_p, k_s) , Σ produces a signature σ . The signing algorithm might be probabilistic.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

CCS'01, November 5-8, 2001, Philadelphia, Pennsylvania, USA.
Copyright 2001 ACM 1-58113-385-5/01/0011 ...\$5.00.

- The *verification algorithm V*. Given a signature σ , a message m and a public key k_p , V tests whether σ is a valid signature of m with respect to k_p . In general, the verification algorithm need not be probabilistic.

2.2 Forgeries and Attacks

In this subsection, we formalize some security notions which capture the main practical situations. On the one hand, the goals of the adversary may be various:

- Disclosing the secret key of the signer. It is the most serious attack. This attack is termed *total break*.
- Constructing an efficient algorithm which is able to sign messages with good probability of success. This is called *universal forgery*.
- Providing a new message-signature pair. This is called *existential forgery*.

In many cases this latter forgery, the *existential forgery*, is not dangerous, because the output message is likely to be meaningless. Nevertheless, a signature scheme which is not existentially unforgeable (and thus that admits existential forgeries) does not guarantee by itself the identity of the signer. For example, it cannot be used to certify randomly looking elements, such as keys. Furthermore, it cannot formally guarantee the non-repudiation property, since anyone may be able to produce a message with a valid signature.

On the other hand, various means can be made available to the adversary, helping her into her forgery. We focus on two specific kinds of attacks against signature schemes: the *no-message attacks* and the *known-message attacks*. In the first scenario, the attacker only knows the public key of the signer. In the second one, the attacker has access to a list of valid message-signature pairs. According to the way this list was created, we usually distinguish many subclasses, but the strongest is the *adaptively chosen-message attack*, where the attacker can ask the signer to sign any message of her choice. She can therefore adapt her queries according to previous answers.

When one designs a signature scheme, one wants to computationally rule out existential forgeries even under adaptively chosen-message attacks, which is the strongest security level for a signature scheme.

3. GENERIC ALGORITHMS

Before we proceed, let us stress that although the generic model in which we analyze DSA offers a somehow weaker form of security than the reductions that we apply to GHR in the standard model, it still provides evidence that twinning may indeed have a beneficial effect on security.

Generic algorithms [15, 21], as introduced by Nechaev and Shoup, encompass group algorithms that do not exploit any special property of the encodings of group elements other than the property that each group element is encoded by a unique string. Typically, algorithms like Pollard's ρ algorithm [18] fall under the scope of this formalism while index-calculus methods do not.

3.1 The Framework

Recall that any Abelian finite group Γ is isomorphic to a product of cyclic groups of the form $(\mathbb{Z}_{p^k}, +)$, where p is a prime. Such groups will be called standard Abelian groups.

An encoding of a standard group Γ is an injective map from Γ into a set of bit-strings S .

We give some examples: consider the multiplicative group of invertible elements modulo some prime q . This group is cyclic and isomorphic to the standard additive group $\Gamma = \mathbb{Z}_{q-1}$. Given a generator g , an encoding σ is obtained by computing the binary representation $\sigma(x)$ of $g^x \bmod q$. The same construction applies when one considers a multiplicative subgroup of prime order r . Similarly, let E be the group of points of some non-singular elliptic curve over a finite field \mathbb{F} , then E is either isomorphic to a (standard) cyclic group Γ or else is isomorphic to a product of two cyclic groups $\mathbb{Z}_{d_1} \times \mathbb{Z}_{d_2}$. In the first case, given a generator G of E , an encoding is obtained by computing $\sigma(x) = x.G$, where $x.G$ denotes the scalar multiplication of G by the integer x and providing coordinates for $\sigma(x)$. The same construction applies when E is replaced by one of its subgroups of prime order r . Note that the encoding set appears much larger than the group size, but compact encodings using only one coordinate and a sign bit ± 1 exist and for such encodings, the image of σ is included in the binary expansions of integers $< tr$ for some small integer t , provided that r is close enough to the size of the underlying field \mathbb{F} . This is exactly what is recommended for cryptographic applications [10].

A generic algorithm A over a standard Abelian group Γ is a probabilistic algorithm that takes as input an encoding list $\{\sigma(x_1), \dots, \sigma(x_k)\}$, where each x_i is in Γ . While it executes, the algorithm may consult an oracle for further encodings. Oracle calls consist of triples $\{i, j, \epsilon\}$, where i and j are indices of the encoding list and ϵ is \pm . The oracle returns the string $\sigma(x_i \pm x_j)$, according to the value of ϵ and this bit-string is appended to the list, unless it was already present. In other words, A cannot access an element of Γ directly but only through its name $\sigma(x)$ and the oracle provides names for the sum or difference of two elements addressed by their respective names. Note however that A may access the list at any time. In many cases, A takes as input a pair $\{\sigma(1), \sigma(x)\}$. Probabilities related to such algorithms are computed with respect to the internal coin tosses of A as well as the random choices of σ and x .

The following theorem appears in [21]:

THEOREM 1. *Let Γ be a standard cyclic group of order N and let p be the largest prime divisor of N . Let A be a generic algorithm over Γ that makes at most n queries to the oracle. If $x \in \Gamma$ and an encoding σ are chosen at random, then the probability that A returns x on input $\{\sigma(1), \sigma(x)\}$ is $O(n^2/p)$.*

PROOF. We refer to [21] for a proof. However, we will need, as an ingredient for our own proofs, the probabilistic model used by Shoup. We develop the model in the special case where N is a prime number r , which is of interest to us. Alternatively, we could work in a subgroup of prime order r .

Basically, we would like to identify the probabilistic space consisting of σ and x with the space $S^{n+2} \times \Gamma$, where S is the set of bit-string encodings. Given a tuple $\{z_1, \dots, z_{n+2}, y\}$ in this space, z_1 and z_2 are used as $\sigma(1)$ and $\sigma(x)$, the successive z_i are used in sequence to answer the oracle queries and the unique value y from Γ serves as x . However, this interpretation may yield inconsistencies as it does not take care of possible collisions between oracle queries. To overcome the difficulty, Shoup defines, along with the execution

of \mathcal{A} , a sequence of linear polynomials $F_i(X)$, with coefficients modulo r . Polynomials F_1 and F_2 are respectively set to $F_1 = 1$ and $F_2 = X$ and the definition of polynomial F_ℓ is related to the ℓ -th query $\{i, j, \epsilon\}$: $F_\ell = F_i \pm F_j$, where the sign \pm is chosen according to ϵ . If F_ℓ is already listed as a previous polynomial F_h , then F_ℓ is marked and \mathcal{A} is fed with the answer of the oracle at the h -th query. Otherwise, z_ℓ is returned by the oracle. Once \mathcal{A} has come to a stop, the value of x is set to y .

It is easy to check that the behavior of the algorithm which plays with the polynomials F_i is exactly similar to the behavior of the regular algorithm, if we require that y is not a root of any polynomial $F_i - F_j$, where i, j range over indices of unmarked polynomials. A sequence $\{z_1, \dots, z_{n+2}, y\}$ for which this requirement is met is called a safe sequence. Shoup shows that, for any $\{z_1, \dots, z_{n+2}\}$, the set of y such that $\{z_1, \dots, z_{n+2}, y\}$ is not safe has probability $\mathcal{O}(n^2/r)$. From a safe sequence, one can define x as y and σ as any encoding which satisfies $\sigma(F_i(y)) = z_i$, for all unmarked F_i . This correspondence preserves probabilities. However, it does not completely cover the sample space $\{\sigma, x\}$ since executions such that $F_i(x) = F_j(x)$, for some indices i, j , such that F_i and F_j are not identical are omitted. To conclude the proof of the above theorem in the special case where N is a prime number r , we simply note that the output of a computation corresponding to a safe sequence $\{z_1, \dots, z_{n+2}, y\}$ does not depend on y . Hence it is equal to y with only minute probability. \square

3.2 Digital Signatures over Generic Groups

We now explain how generic algorithms can deal with attacks against DSA-like signature schemes [6, 20, 16, 10]. We do this by defining a generic version of DSA that we call GDSA. Parameters for the signature include a standard cyclic group of prime order r together with an encoding σ . The signer also uses as a secret key/public key pair $\{x, \sigma(x)\}$. Note that we have chosen to describe signature generation as a regular rather than generic algorithm, using a full description of σ . To sign a message m , $1 < m < r$ the algorithm executes the following steps:

1. Generate a random number u , $1 \leq u < r$.
2. Compute $c \leftarrow \sigma(u) \bmod r$. If $c = 0$ go to step 1.
3. Compute $d \leftarrow u^{-1}(m + xc) \bmod r$. If $d = 0$ go to step 1.
4. Output the pair $\{c, d\}$ as the signature of m .

The verifier, on the other hand, is generic:

1. If $c \notin [1, r-1]$ or $d \notin [1, r-1]$, output invalid and stop.
2. Compute $h \leftarrow d^{-1} \bmod r$, $h_1 \leftarrow hm \bmod r$ and $h_2 \leftarrow hc \bmod r$.
3. Obtain $\sigma(h_1 + h_2x)$ from the oracle and compute $c' \leftarrow \sigma(h_1 + h_2x) \bmod r$.
4. If $c \neq c'$ output invalid and stop otherwise output valid and stop.

The reader may wonder how the verifier obtains the value of σ requested at step 3. This is simply achieved by mimicking the usual double-and-add algorithm and asking the

appropriate queries to the oracle. This yields $\sigma(h_1)$ and $\sigma(h_2x)$. A final call to the oracle completes the task.

A generic algorithm \mathcal{A} can also perform forgery attacks against a signature scheme. This is defined by the ability of \mathcal{A} to return on input $\{\sigma(1), \sigma(x)\}$ a triple $\{m, c, d\} \in \Gamma^3$ for which the verifier outputs valid. Here we assume that both algorithms are performed at a stretch, keeping the same encoding list.

To deal with adaptive attacks one endows \mathcal{A} with another oracle, called the signing oracle. To query this oracle, the algorithm provides an element $m \in \Gamma$. The signing oracle returns a valid signature $\{c, d\}$ of m . Success of \mathcal{A} is defined by its ability to produce a valid triple $\{\tilde{m}, \tilde{c}, \tilde{d}\}$, such that \tilde{m} has not been queried during the attack.

Such a forgery can be easily performed against this GDSA scheme, even with just a passive attack: the adversary chooses random numbers h_1 and h_2 , $1 \leq h_1, h_2 < r$ and computes $c \leftarrow \sigma(h_1 + h_2x) \bmod r$. Then it defines $d = ch_2^{-1} \bmod r$, $h = d^{-1} \bmod r$, and eventually $m = dh_1 \bmod r$. The triple $\{m, c, d\} \in \Gamma^3$ is therefore a valid one, unless $c = 0$, which is very unlikely.

4. THE SECURITY OF TWIN GDSA

4.1 A Theoretical Result

The above definitions extend to the case of twin signatures, by requesting the attacker \mathcal{A} to output an m and two distinct pairs $\{c, d\} \in \Gamma^2$, $\{c', d'\} \in \Gamma^2$. Success is granted as soon as the verifying algorithm outputs valid for both triples¹. We prove the following:

THEOREM 2. *Let Γ be a standard cyclic group of prime order r . Let S be a set of bit-string encodings of cardinality at least r , included in the set of binary representations of integers $< tr$, for some t . Let \mathcal{A} be a generic algorithm over Γ that makes at most n queries to the oracle. If $x \in \Gamma$ and an encoding σ are chosen at random, then the probability that \mathcal{A} returns a message m together with two distinct GDSA signatures of m on input $\{\sigma(1), \sigma(x)\}$ is $\mathcal{O}(tn^2/r)$.*

PROOF. We cover the non adaptive case and tackle the more general case after the proof. We use the probabilistic model developed in section 3.1. Let \mathcal{A} be a generic attacker able to forge some m and two distinct signatures $\{c, d\}$ and $\{c', d'\}$. We assume that, once these outputs have been produced, \mathcal{A} goes on checking both signatures; we estimate the probability that both are valid.

We restrict our attention to behaviors of the full algorithm corresponding to safe sequences $\{z_1, \dots, z_{n+2}, y\}$. By this, we discard a set of executions of probability $\mathcal{O}(n^2/r)$. We let P be the polynomial $(md^{-1}) + (cd^{-1})X$ and Q be the polynomial $(md'^{-1}) + (c'd'^{-1})X$.

- We first consider the case where either P or Q does not appear in the F_i list before the signatures are produced. If this happens for P , then P is included in the F_i list at signature verification and the corresponding answer of the oracle is a random number z_i . Unless $z_i = c \bmod r$, which is true with probability at most

¹ using [14] the simultaneous square-and-multiply generation or verification of two DSA signatures is only 17% slower than the generation or verification of a single signature.

t/r , the signature is invalid. A similar bound holds for Q .

- We now assume that both P and Q appear in the F_i list before \mathcal{A} outputs its signatures. We let i denote the first index such that $F_i = P$ and j the first index such that $F_j = Q$. Note that both F_i and F_j are unmarked (as defined in section 3.1). If $i = j$, then we obtain that $md^{-1} = md'^{-1}$ and $cd^{-1} = c'd'^{-1}$. From this, it follows that $c = c'$, $d = d'$ and the signatures are not distinct.
- We are left with the case where $i \neq j$. We let $\Omega_{i,j}$, $i < j$, be the set of safe sequences producing two signatures such that the polynomials P , Q , defined as above appear for the first time before the algorithm outputs the signatures, as F_i and F_j . We consider a fixed value w for $\{z_1, \dots, z_{j-1}\}$ and let \hat{w} be the set of safe sequences extending w . We note that F_i and F_j are defined from w and we write $F_i = a + bX$, $F_j = a' + b'X$. We claim that $\Omega_{i,j} \cap \hat{w}$ has probability $\leq t/r$. To show this, observe that one of the signatures that the algorithm outputs is necessarily of the form $\{c, d\}$, with $c = z_i \bmod r$, $c = db \bmod r$ and $m = da \bmod r$. Now, the other signature is $\{c', d'\}$ and since m is already defined we get $d' = ma'^{-1} \bmod r$ and $c' = b'd' \bmod r$. This in turn defines $z_j \bmod r$ within a subset of at most t elements. From this, the required bound follows and, from the bound, we infer that the probability of $\Omega_{i,j}$ is at most t/r .

Summing up, we have bounded the probability that a safe sequence produces an execution of \mathcal{A} outputting two valid signatures by $\mathcal{O}(tn^2/r)$. This finishes the proof. \square

In the proof, we considered the case of an attacker forging a message-signature pair from scratch. A more elaborate scenario corresponds to an attacker who can adaptively request twin signatures corresponding to messages of his choice. In other words, the attacker interacts with the legitimate signer by submitting messages selected by its program.

We show how to modify the security proof that was just given to cover the adaptive case. We assume that each time it requests a signature the attacker \mathcal{A} immediately verifies the received signature. We also assume that the verification algorithm is normalized in such a way that, when verifying a signature $\{c, d\}$ of a message m , it asks for $\sigma((md^{-1}) + (cd^{-1})x)$ after a fixed number of queries, say q . We now explain how to simulate signature generation: as before, we restrict our attention to behaviors of the algorithm corresponding to safe sequences $\{z_1, \dots, z_{n+2}, y\}$. When the (twin) signature of m is requested at a time of the computation when the encoding list contains i elements, one picks z_{i+q} and z_{i+2q} and manufactures the two signatures as follows:

1. Let $c \leftarrow z_{i+q} \bmod r$, pick d at random.
2. Let $c' \leftarrow z_{i+2q} \bmod r$, pick d' at random.
3. Output $\{c, d\}$ and $\{c', d'\}$ as the first and second signatures.

While verifying both signatures, \mathcal{A} will receive the elements z_{i+q} and z_{i+2q} , as

$$\sigma((md^{-1}) + (cd^{-1})x) \text{ and } \sigma((md'^{-1}) + (c'd'^{-1})x)$$

respectively, unless F_{i+q} or F_{i+2q} appears earlier in the F_i list. Due to the randomness of d and d' , this happens with very small probability bounded by n/r . Altogether, the simulation is spotted with probability $\mathcal{O}(n^2/r)$ which does not affect the $\mathcal{O}(tn^2/r)$ bound for the probability of successful forgery.

4.2 Practical Meaning of the Result

We have shown that, in the setting of generic algorithms, existential forgery against twin GDSA has a minute success probability. Of course this does not tell anything on the security of actual twin DSA. Still, we believe that our proof has some practical meaning. The analogy with hash functions and the random oracle model [3] is inspiring: researchers and practitioners are aware that proofs in the random oracle model are not proofs but a mean to spot design flaws and validate schemes that are supported by such proofs. Still, all standard signature schemes that have been proposed use specific functions which are not random by definition; our proofs seem to indicate that if existential forgery against twin DSA is possible, it will require to dig into structural properties of the encoding function. This is of some help for the design of actual schemes: for example, the twin DSA described in Appendix A allows signature with message recovery without hashing and without any form of redundancy, while keeping some form of provable security. This might be considered a more attractive approach than [17] or [1], the former being based on redundancy and the latter on random oracles. We believe that twin DSA is even more convincing in the setting of elliptic curves, where there are no known ways of taking any advantage of the encoding function.

5. AN RSA-BASED TWINNING IN THE STANDARD MODEL

The twin signature scheme described in this section belongs to the (very) short list of efficient schemes provably secure in the standard model: in the sequel, we show that producing existential forgeries even under an adaptively chosen-message attack is equivalent to solving the Flexible RSA Problem [5].

Security in the standard model implies no ideal assumptions; in other words we directly reduce the Flexible RSA Problem to a forgery. As a corollary, we present an efficient and provably secure signature scheme that does not require any hash function.

Furthermore, the symmetry provided by twinning is much simpler to analyze than Cramer-Shoup's proposal [5] which achieves a similar security level, and similar efficiency, with a rather intricate proof.

5.1 Gennaro-Halevi-Rabin Signatures

In [8] Gennaro, Halevi and Rabin present the following signature scheme: Let n be an ℓ -bit RSA modulus [19], H a hash-function and $y \in \mathbb{Z}_n^*$. The pair $\{n, y\}$ is the signer's public key, whose secret key is the factorization of n .

- To sign m , the signer hashes $e \leftarrow H(m)$ (which is very likely to be co-prime with $\varphi(n)$) and computes the e -th

root of y modulo n using the factorization of n :

$$s \leftarrow y^{1/e} \bmod n$$

- To verify a given $\{m, s\}$, the verifier checks that

$$s^{H(m)} \bmod n \stackrel{?}{=} y.$$

Security relies on the Strong RSA Assumption. Indeed, if H outputs elements that contain at least a new prime factor, existential forgery is impossible. Accordingly, Genaro et al. define a new property that H must satisfy to yield secure signatures: *division intractability*. Division intractability means that it is computationally impossible to find a_1, \dots, a_k and b such that $H(b)$ divides the product of all the $H(a_i)$. In [8], it is conjectured that such functions exist and heuristic conversions from collision-resistant into division-intractable functions are shown (see also [4]).

Still, security against adaptively chosen-message attacks requires the hash function H to either behave like a random oracle model or achieve the chameleon property [12]. This latter property, for a hash function, provides a trapdoor which helps to find second preimages, even with some fixed part. Indeed, some signatures can be pre-computed, but with specific exponents before outputting $y: y = x^{\prod e_i} \bmod n$ for random primes $e_i = H(m_i, r_i)$.

Using the chameleon property, for the i -th query m to the signing oracle, the simulator who knows the trapdoor can get an r such that $H(m_i, r_i) = H(m, r) = e_i$. In the random oracle model, one simply defines $H(m, r) \leftarrow e_i$.

Then $s = x^{\prod_{i=1}^k e_i} = y^{1/e_i} \bmod n$ and the signature therefore consists of the triple $\{m, r, s\}$ satisfying

$$s^{H(m,r)} = y \bmod n.$$

Cramer and Shoup [5] also proposed a scheme based on the Strong RSA Assumption, the first practical signature scheme to be secure in the standard model, but with universal one-way hash functions; our twin scheme will be similar but with a nice symmetry in the description (which helps for the security analysis) and no hash-functions, unless one wants to sign a long message.

5.2 Preliminaries

We build our scheme in two steps. The first scheme resists existential forgeries when subjected to no-message attacks. Twinning will immune it against adaptively chosen-message attacks.

5.2.1 Injective function into the prime integers.

Before any description, we will assume the existence of a function p with the following properties: given a security parameter k (which will be the size of the signed messages), p maps any string from $\{0, 1\}^k$ into the set of the prime integers, p is also designed to be easy to compute and injective. A candidate is proposed and analyzed in Appendix B.

5.2.2 The Flexible RSA Problem and the Strong RSA Assumption.

Let us also recall the *Flexible RSA Problem* [5]. Given an RSA modulus n and an element $y \in \mathbb{Z}_n^*$, find any exponent $e > 1$, together with an element x such that $x^e = y \bmod n$.

The Strong RSA Assumption is the conjecture that this problem is intractable for large moduli. This was indepen-

dently introduced by [2, 7], and then used in many further security analyses (e.g. [5, 8]).

5.3 A First GHR Variant

The first scheme is very similar to GHR without random oracles but with function p instead:

- To sign $m \in \{0, 1\}^k$, the signer computes $e \leftarrow p(m)$ and the e -th root of y modulo n using the factorization of n

$$s \leftarrow y^{1/e} \bmod n$$

- To verify a given $\{m, s\}$, the verifier checks that

$$s^{p(m)} \bmod n \stackrel{?}{=} y.$$

Since p provides a new prime for each new message (injectivity), existential forgery contradicts the Strong RSA Assumption. However, how can we deal with adaptively chosen-message attacks without any control over the output of the function p , which is a publicly defined non-random oracle and not a trapdoor function either?

5.4 The Twin Version

The final scheme is quite simple since it consists in duplicating the previous one: the signer uses two ℓ -bit RSA moduli n_1, n_2 and two elements y_1, y_2 in $\mathbb{Z}_{n_1}^*$ and $\mathbb{Z}_{n_2}^*$ respectively. Secret keys are the prime factors of the n_i .

- To sign a message m , the signer probabilistically derives two messages $\mu_1, \mu_2 \in \{0, 1\}^k$, (from m and a random tape ω), computes $e_i \leftarrow p(\mu_i)$ and then the e_i -th root of y_i modulo n_i , for $i = 1, 2$, using the factorization of the moduli:

$$\{s_1 \leftarrow y_1^{1/e_1} \bmod n_1, s_2 \leftarrow y_2^{1/e_2} \bmod n_2\}$$

- To verify a given $\{m, \omega, s_1, s_2\}$, the verifier computes μ_1 and μ_2 , then checks that $s_i^{p(\mu_i)} \bmod n_i \stackrel{?}{=} y_i$, for $i = 1, 2$.

To prevent forgeries, a new message must involve a new exponent, either e_1 or e_2 , which never occurred in the signatures provided by the signing oracle. Therefore, a first requirement is that μ_1 and μ_2 define at most one message m , but only if they have been correctly constructed. Thus, some redundancy is furthermore required.

We thus suggest the following derivation, to get μ_1 and μ_2 from $m \in \{0, 1\}^{k/2}$ (we assume k to be even): one chooses two random elements $a, b \in \{0, 1\}^{k/2}$, then $\mu_1 = (m \oplus a) \parallel (m \oplus b)$ and $\mu_2 = a \parallel b$.

Clearly, given μ_1 and μ_2 , one gets back $M = \mu_1 \oplus \mu_2$, which provides a valid message if and only if the redundancy holds: $\bar{M} = \underline{M}$, where \bar{S} and \underline{S} denote the two $k/2$ -bit halves of a k -bit string S , the most significant and the least significant parts respectively.

5.5 Existential Forgeries

Let us show that existential forgery of the twin scheme, with above derivation process, leads to a new solution of the Flexible RSA Problem:

LEMMA 1. After q queries to the signing oracle, the probability that there exist a new message m and values a, b ,

which lead to $\mu_1 = (m \oplus a) \parallel (m \oplus b)$ and $\mu_2 = a \parallel b$, such that both $e_1 = p(\mu_1)$ and $e_2 = p(\mu_2)$ already occurred in the signatures provided by the signing oracle is less than $q^2/2^{k/2}$.

PROOF. Let $\{m_i, a_i, b_i, s_{1,i}, s_{2,i}\}$ denote the answers of the signing oracle. Using the injectivity of p , the existence of such m , a and b means that there exist indices i and j for which

$$\begin{aligned} (m \oplus a) \parallel (m \oplus b) = \mu_1 &= \mu_{1,i} = (m_i \oplus a_i) \parallel (m_i \oplus b_i) \\ a \parallel b = \mu_2 &= \mu_{2,j} = a_j \parallel b_j. \end{aligned}$$

Then

$$a \oplus b = (m \oplus a) \oplus (m \oplus b) = (m_i \oplus a_i) \oplus (m_i \oplus b_i) = a_i \oplus b_i, \text{ and}$$

$$a \oplus b = a_j \oplus b_j.$$

Therefore, for a $j > i$ (the case $i > j$ is similar), the new random elements a_j, b_j must satisfy $a_j \oplus b_j = a_i \oplus b_i$. Since it is randomly chosen by the signer, the probability that this occurs for some $i < j$ is less than $(j-1)/2^{k/2}$.

Altogether, the probability that for some j there exists some $i < j$ which satisfies the above equality is less than $q^2/2 \times 2^{-k/2}$. By symmetry, we obtain the same result if we exchange i and j .

The probability that both exponents already appeared is consequently smaller than $q^2/2^{k/2}$. \square

To prevent adaptively chosen-message attacks, we need no trapdoor property for p , nor random oracle assumption either. We simply give the factorization of one modulus to the simulator, which can use any pre-computed exponentiation with any new message, as when chameleon functions are used [8].

5.6 Adaptively Chosen-Message Attacks

Indeed, to prevent adaptively chosen-message attacks, one just needs to describe a simulator; our simulator works as follows:

- The simulator is first given the moduli n_1, n_2 and the elements $y_1 \in \mathbb{Z}_{n_1}^*$, $y_2 \in \mathbb{Z}_{n_2}^*$, as well as the factorization of n_1 , where γ is randomly chosen in $\{1, 2\}$. To simplify notations we assume that $\gamma = 1$. And the following works without loss of generality since the derivation of μ_1 and μ_2 is perfectly symmetric: they are randomly distributed, but satisfy $\mu_1 \oplus \mu_2 = m \parallel m$ (it is a perfect secret sharing).
- The simulator randomly generates q values $e_{2,j} \leftarrow p(\mu_{2,j})$, with randomly chosen $\mu_{2,j} \in_R \{0, 1\}^k$ for $j = 1, \dots, q$ and computes

$$z \leftarrow y_2^{\prod_{j=1, \dots, q} e_{2,j}} \bmod n_2.$$

The new public key for the signature scheme is the following: the moduli n_1, n_2 with the elements y_1, z in $\mathbb{Z}_{n_1}^*$ and $\mathbb{Z}_{n_2}^*$ respectively.

- For the j -th signed message m , the simulator first gets $(a \parallel b) \leftarrow (m \parallel m) \oplus \mu_{2,j}$. It therefore computes $\mu_1 \leftarrow a \parallel b$, and thus $\mu_2 \leftarrow \mu_{2,j} = (m \oplus a) \parallel (m \oplus b)$.

Then, it knows $s_2 = y_2^{\prod_{i \neq j} e_{2,i}} \bmod n_2$, and computes s_1 using the factorization of n_1 .

Such a simulator can simulate up to q signatures, which leads to the following theorem.

THEOREM 3. *Let us consider an adversary against the twin-GHR scheme who succeeds in producing an existential forgery, with probability greater than ϵ , after q adaptive queries to the signing oracle in time t , then the Flexible RSA Problem can be solved with probability greater than ϵ' within a time bound t' , where*

$$\epsilon' = \frac{1}{2} \left(\epsilon - \frac{q^2}{2^{k/2}} \right) \quad \text{and} \quad t' = t + O(q \times \ell^2 \times k).$$

PROOF. Note that the above bounds are almost optimal since $\epsilon' \cong \epsilon/2$ and $t' \cong 2t$. Indeed, the time needed to produce an existential forgery after q signature queries is already in $O(q \times (|n_1|^2 + |n_2|^2)k)$. To evaluate the success probability, q is less than say 2^{40} , but k may be taken greater than 160 bits (and even much more).

To conclude the proof, one just needs to address the random choice of γ . As we have seen in Lemma 1, with probability greater than $\epsilon - q^2/2^{k/2}$, one of the exponents in the forgery never appeared before. Since γ is randomly chosen and the view of the simulation is perfectly independent of this choice, with probability of one half, $e = e_1$ is new. Let us follow our assumption that $\gamma = 1$, then

$$s^e = s_2^e = z = y_2^\pi \bmod n_2,$$

where $\pi = \prod_{j=1, \dots, q} e_{2,j}$. Since e is new, it is relatively prime with π , and therefore, there exist u and v such that $ue + v\pi = 1$: let us define $x = y_2^u s^v \bmod n_2$,

$$x^e = (y_2^u s^v)^e = y_2^{1-v\pi} s^{ev} = y_2 (y_2^\pi)^{-v} (s^e)^v = y_2 \bmod n_2.$$

We thus obtain an e -th root of the given y_2 modulo n_2 , for a new prime e . \square

5.7 More Signatures

One may remark that the length of the messages we can sign with above construction is limited to $k/2$ bits, because of the required redundancy. But one can increase the size, by signing three derived messages: in order to sign $m \in \{0, 1\}^k$, one chooses two random elements $a, b \in \{0, 1\}^{k/2}$ (we still assume k to be even), and signs with different moduli

$$\begin{aligned} \mu_1 &= m \oplus (a \parallel b) \\ \mu_2 &= a \parallel b \\ \mu_3 &= m \oplus (b \parallel a). \end{aligned}$$

6. CONCLUSION

AND FURTHER RESEARCH

We proposed an alternative to the well-known hash-and-sign paradigm, based on the simple idea of signing twice (or more) identical or related short messages. We believe that our first investigations show that this is a promising strategy, deserving further study.

A number of interesting questions remain open. First, from the efficiency point of view, which is a frequent concern, we are aware that the current proposals do not deal with either the computational cost, or the communication load, in an efficient way. Thus, for example, can the number of fields in a twin DSA be reduced from four ($\{c, d\}$ and $\{c', d'\}$) to three or less? Can we also suppress some fields in the twin-GHR, or sign k -bit long messages with only two signatures?

Finally, can an increase in the number of signatures (e.g. three instead of two) yield better security bounds?

7. REFERENCES

- [1] M. Abe and T. Okamoto. A Signature Scheme with Message Recovery as Secure as Discrete Logarithm. In *Asiacrypt '99*, LNCS 1716. Springer-Verlag, Berlin, 1999.
- [2] N. Barić and B. Pfitzmann. Collision-Free Accumulators and Fail-Stop Signature Schemes without Trees. In *Eurocrypt '97*, LNCS 1233, pages 480–484. Springer-Verlag, Berlin, 1997.
- [3] M. Bellare and P. Rogaway. Random Oracles Are Practical: a Paradigm for Designing Efficient Protocols. In *Proc. of the 1st CCS*, pages 62–73. ACM Press, New York, 1993.
- [4] J.-S. Coron and D. Naccache. Security Analysis of the Gennaro-Halevi-Rabin Signature Scheme. In *Eurocrypt '99*, LNCS 1592, pages 91–101. Springer-Verlag, Berlin, 1999.
- [5] R. Cramer and V. Shoup. Signature Scheme based on the Strong RSA Assumption. In *Proc. of the 6th CCS*, pages 46–51. ACM Press, New York, 1999.
- [6] T. El Gamal. A Public Key Cryptosystem and a Signature Scheme Based on Discrete Logarithms. *IEEE Transactions on Information Theory*, IT-31(4):469–472, July 1985.
- [7] E. Fujisaki and T. Okamoto. Statistical Zero Knowledge Protocols to Prove Modular Polynomial Relations. In *Crypto '97*, LNCS 1294, pages 16–30. Springer-Verlag, Berlin, 1997.
- [8] R. Gennaro, S. Halevi, and T. Rabin. Secure Hash-and-Sign Signature Without the Random Oracle. In *Eurocrypt '99*, LNCS 1592, pages 123–139. Springer-Verlag, Berlin, 1999.
- [9] S. Goldwasser, S. Micali, and R. Rivest. A Digital Signature Scheme Secure Against Adaptive Chosen-Message Attacks. *SIAM Journal of Computing*, 17(2):281–308, April 1988.
- [10] IEEE P1363. Standard Specifications for Public Key Cryptography. Available from <http://grouper.ieee.org/groups/1363>, August 1998.
- [11] M. Jakobsson and C. P. Schnorr. Security of Discrete Logarithm Cryptosystems in the Random Oracle Model and Generic Model. Available from <http://www.bell-labs.com/~markusj>, 1998.
- [12] H. Krawczyk and T. Rabin. Chameleon Hashing and Signatures. In *Proc. of NDSS '2000*. Internet Society, 2000.
- [13] A. Menezes, P. van Oorschot, and S. Vanstone. *Handbook of Applied Cryptography*. CRC Press, 1996. Available from <http://www.cacr.math.uwaterloo.ca/hac/>.
- [14] D. M'Raihi and D. Naccache. Batch Exponentiation – A Fast DLP-based Signature Generation Strategy. In *Proc. of the 3rd CCS*, pages 58–61. ACM Press, New York, 1996.
- [15] V. I. Nechaev. Complexity of a Determinate Algorithm for the Discrete Logarithm. *Mathematical Notes*, 55(2):165–172, 1994.
- [16] NIST. Digital Signature Standard (DSS). Federal Information Processing Standards Publication 186, November 1994.
- [17] K. Nyberg and R. A. Rueppel. Message Recovery for Signature Schemes Based on the Discrete Logarithm Problem. In *Eurocrypt '94*, LNCS 950, pages 182–193. Springer-Verlag, Berlin, 1995.
- [18] J. M. Pollard. Monte Carlo Methods for Index Computation (mod p). *Mathematics of Computation*, 32(143):918–924, July 1978.
- [19] R. Rivest, A. Shamir, and L. Adleman. A Method for Obtaining Digital Signatures and Public Key Cryptosystems. *Communications of the ACM*, 21(2):120–126, February 1978.
- [20] C. P. Schnorr. Efficient Signature Generation by Smart Cards. *Journal of Cryptology*, 4(3):161–174, 1991.
- [21] V. Shoup. Lower Bounds for Discrete Logarithms and Related Problems. In *Eurocrypt '97*, LNCS 1233, pages 256–266. Springer-Verlag, Berlin, 1997.

APPENDIX

A. TWIN SIGNATURES WITH MESSAGE RECOVERY

In this appendix, we describe a twin version of the Nyberg-Rueppel scheme [17] which provides message recovery. Keeping the notations of section 4.1:

1. Generate a random number u , $1 \leq u < r$.
2. Compute $c \leftarrow \sigma(u) + m \bmod r$. If $c = 0$ go to step 1.
3. Compute an integer $d \leftarrow u - cx \bmod r$.
4. Output the pair $\{c, d\}$ as the signature.

In the above, f is what is called in [10] a message with appendix. It simply means that it has an adequate redundancy. The corresponding verification is performed by the following (generic) steps:

1. If $c \notin [1, r-1]$ or $d \notin [0, r-1]$, output invalid and stop.
2. Obtain $\sigma(d + cx)$ from the oracle and compute $\gamma \leftarrow \sigma(d + cx) \bmod r$.
3. Check the redundancy of $m \leftarrow c - \gamma \bmod r$. If incorrect output invalid and stop; otherwise output the reconstructed message m , output valid and stop.

In the twin setting, signature generation is alike but is performed twice, so as to output two distinct signatures. However, no redundancy is needed. The verifier simply checks that the signatures are distinct and outputs two successive versions of the message, say m and m' . It returns valid if $m \neq m'$ and invalid otherwise. The security proof is sketched here, we leave the discussion of adaptive attacks to the reader.

We keep the notations and assumptions of section 4 and let \mathcal{A} be a generic attacker over Γ which outputs, on input $\{\sigma(1), \sigma(x)\}$, two signature pairs $\{c, d\}$, $\{c', d'\}$ and runs the verifying algorithm that produces from these signatures two messages m , m' and checks whether they are equal. We wish to show that, if $x \in \Gamma$ and an encoding σ are chosen at random, then the probability that $m = m'$ is $O((n^2/r))$.

As before, we restrict our attention to behaviors of the full algorithm corresponding to safe sequences $\{z_1, \dots, z_n, y\}$.

We let P, Q be the polynomials $d + cX$ and $d' + c'X$. We first consider the case where either P or Q does not appear in the F_i list before the signatures are produced. If this happens for P , then, P is included in the F_i list at signature verification and the corresponding answer of the oracle is a random number z_i . Since m is computed as $c - z_i \bmod r$, the probability that $m = m'$ is bounded by t/r . A similar bound holds for Q .

We now assume that both P and Q appear in the F_i list before A outputs its signatures. We let i denote the first index such that $F_i = P$ and j the first index such that $F_j = Q$. Note that both F_i and F_j are unmarked (as defined in section 3.1). If $i = j$, then we obtain that $c = c'$ and $d = d'$. From this, it follows that the signatures are not distinct.

As in section 4, we are left with the case where $i \neq j$ and we define $\Omega_{i,j}$, $i < j$, to be the set of safe sequences producing two signatures such that the polynomials P, Q , defined as above appear for the first time before the algorithm outputs the signatures, as F_i and F_j . We show that, for any fixed value $w = \{z_1, \dots, z_{j-1}\}$, $\Omega_{i,j} \cap \bar{w}$ has probability $\leq t/r$, where \bar{w} is defined as above. Since we have $m = c - z_i \bmod r$ and $m' = c' - z_j \bmod r$, we obtain $z_j = c' - c + z_i \bmod r$, from which the upper bound follows. From this bound, we obtain that the probability of $\Omega_{i,j}$ is at most t/r and, taking the union of the various $\Omega_{i,j}$ s, we conclude that the probability to obtain a valid twin signature is at most $O(tn^2/r)$.

B. THE CHOICE OF FUNCTION P

B.1 A Candidate

The following is a natural candidate:

$$p: \{0, 1\}^k \rightarrow \mathbb{P} \\ m \mapsto \text{nextprime}(m \times 2^r)$$

where r is suitably chosen to guarantee the existence of a prime in any set $[m \times 2^r, (m+1) \times 2^r]$, for $m < 2^k$.

Note that the deterministic property of nextprime is not mandatory, one just needs it to be injective. But then, the preimage must be easily recoverable from the prime: the exponent is sent as the signature, from which one checks the primality and extracts the message (message-recovery).

B.2 Analysis

It is clear that any generator of random primes, using m as a seed, can be considered as a candidate for p . The function proposed above is derived from a technique for accelerating prime generation called *incremental search* (e.g. [13], page 148).

1. Input: an odd k -bit number n_0 (derived from m)
2. Test the s numbers $n_0, n_0 + 2, \dots, n_0 + 2(s-1)$ for primality

Under reasonable number-theoretic assumptions, if $s = c \cdot \ln 2^k$, the probability of failure of this technique is smaller than $2e^{-2c}$, for large k .

Using our notations, in such a way that there exists at least a prime in any set $[m \times 2^r, (m+1) \times 2^r]$, but with probability smaller than 2^{-80} , we obtain from above formulae that $c \cong 40$, and $2^r \geq 40 \ln 2^{k+r+1}$. Therefore, a suitable

candidate is $r \cong 5 \log_2 k$, and less than $20k$ primality tests have to be performed.

B.3 Extensions

B.3.1 Collision-resistance:

To sign large messages (at the cost of extra assumptions), one can of course use any collision-resistant hash-function h before signing (using the classical hash-and-sign technique). Clearly, the new function $m \mapsto p(h(m))$ is not mathematically injective, but just computationally injective (which is equivalent to collision-resistance), which is enough for the proof.

B.3.2 Division intractability:

If one wants to improve efficiency, using the division-intractability conjecture proposed in [8], any function that outputs k -bit strings can be used instead of p . More precisely:

Definition (Division Intractability). A function H is said (n, ν, τ) -division intractable if any adversary which runs in time τ cannot find, with probability greater than ν , a set of elements a_1, \dots, a_n and b such that $H(b)$ divides the product of all the $H(a_i)$.

As above, that function p would not be injective, but just collision-resistant, which is enough to prove the following:

THEOREM 4. Let us consider the twin-GHR scheme where p is any (q, ϵ, t) -division-intractable hash function. Let us assume that an adversary A succeeds in producing an existential forgery under an adaptively chosen-message attack within time t and with probability greater than ϵ , after q queries to the signing oracle. Then one can either contradict the division-intractability assumption or solve the Flexible RSA Problem with probability greater than ϵ' within a time bound t' , where

$$\epsilon' = \frac{1}{2} \left(\epsilon - \frac{q^2}{2^{k/2}} \right) \quad \text{and} \quad t' = t + O(q \times \ell^2 \times k).$$

Batch Exponentiation - A Fast DLP-based signature generation strategy -

David M'RAÏHI

David NACCACHE

GEMPLUS, Crypto Team, 1 place de Méditerranée
F-95208, Sarcelles CEDEX, FRANCE
[100145.2261 and 100142.3240]@compuserve.com

Abstract : The signature generation phase of most DLP-based signature schemes (for instance Schnorr[10], El-Gamal[4] or the newly standardized D.S.A.[3]) includes the time-consuming computation of $r = g^k \bmod p$ where k is random.

This paper introduces a new computational strategy that can apply in this particular context :

A *batch exponentiation* technique which allows the generation of large sets of exponentials without introducing any bias between the k s (that is, the signer can batch-compute the exponentials corresponding to arbitrarily imposed powers -for instance by an external random number generator). Our method offers real improvements over the prior art with various time and memory trade-offs.

1. Introduction

In many DLP-based signature schemes¹ the signer performs the operation $r = g^k \bmod p$ where k is random. As the signer is often the "weak party" in the signature protocol, several authors tried to accelerate the exponentiation by pre-computing values [1], [5] or sub-contracting a part of the exponentiation workload to the verifier [6] (provided that a set of precautions is taken into consideration). Except the fact that some of these algorithms were broken [7], [8], extra memory storage is frequently an unrealistic assumption.

In this paper, we investigate a strategy for improving the generation of r : the method (providing improvements ranging from 42% to 85% over the *square-&-multiply* algorithm) can apply to the batch generation of fixed- g -based signatures without introducing any bias into the exponents (that is, we assume that the k s are imposed to the signer by some random source). We assume that no pre-computation is allowed other than what needed to execute similar size basic *square-&-multiply*. The new method may as well open the way to interesting developments for accelerating the computation of discrete logarithms.

Permission to make digital/hard copies of all or part of this material for personal or classroom use is granted without fee provided that the copies are not made or distributed for profit or commercial advantage, the copyright notice, the title of the publication and its date appear, and notice is given that copyright is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires specific permission and/or fee.

CCS '96, New Delhi, India
© 1996 ACM 0-89791-829-0/96/03...\$3.50

2. Unbiased Batch-Exponentiation

The simultaneous signature of many messages by a shop terminal or the processing of electronic documents by an administration frequently involves massive computations consisting in the repetition of small operations. The re-combination of these computations for minimizing the computational effort is thus an interesting research direction (see for instance [10]).

The batch exponentiation technique proposed hereafter is built around the following observation : since the exponents are random, the uniform distribution of ones in different exponents is expected to have some matching patterns. Considering this fact, our batch-exponentiation strategy consists in minimizing the signer's workload by exponentiating the intersection separately and resetting the corresponding bits in the initial exponents.

2.1 Parallel square-&-multiply (straightforward)

Let n be the size of the exponents and N the number of exponentials to be computed. The usual method to generate the r s consists in calculating successive squares of g and performing the required multiplications selected by the bits of each k_i (about $n/2$ multiplications).

Let $S = \{k_i | i=1, \dots, N\}$ be a set of random powers. The corresponding set of exponentials

$$R = \{r_i | i=1, \dots, N\}$$

can be computed by performing the successive squares of g only once. Thus, the total computational effort mainly relies on the average number of multiplications, depending on the Hamming weight of each random exponent.

The algorithm $\text{PSM}(S, R)$, of complexity

$$E(N) = N(n/2 - 1) + (n - 1)$$

using $N + 1$ registers, is :

```
for  $i \leftarrow 1$  to  $N$     $r_i \leftarrow 1$ 
for  $j \leftarrow 0$  to  $n-1$ 
for  $i \leftarrow 1$  to  $N$    if  $k_i[j] = 1$  then  $r_i \leftarrow (r_i * g) \bmod p$ 
 $g \leftarrow g^2 \bmod p$ 
```

$$\text{where } k_i = \sum_{j=0}^{n-1} k_i[j] 2^j.$$

2.2 The Basic Strategy

Denoting S a set of N random powers k , the strategy to generate the related exponentials is the following :

① Cut S into $T = \lfloor N / L \rfloor$ sets $s_h = \{k_j\}_{1 \leq j \leq L}$ where $L \leq 5^2$

② Let $P(s_h) = \bigcup_{i=1}^{2L-1} s_{h,i}$ where $s_{h,i}$ are s_h 's subsets with $s_{h,i} \neq \emptyset$

¹ for instance [3], [4] and [10]

² L value depends on the exponent length; further explanation will be found in 3.2

• For $h \leftarrow 1$ to T

$$\textcircled{D} \text{ For } i \leftarrow 1 \text{ to } 2^L - 1 \quad c_{h,i} = \bigoplus_{k_j \in s_{h,i}} k_j$$

• For $C \leftarrow L-1$ to 1

For $i \leftarrow 1$ to $2^L - 1$

if $\text{Card}(s_{h,i}) = C$ then

$$c_{h,i} = c_{h,i} - \sum_{j \neq i, s_{h,i} \subset s_{h,j}} c_{h,j}$$

• $\text{PSM}(\{c_{h,i}\}, R)$

• For $i \leftarrow 1$ to L

$$g^{k_i} = \prod_{k_j \in s_{h,i}} r_j \quad \text{where } r_j \in R$$

The idea is to operate on each subset, resetting the common bits of the powers, and compute the exponentials together to save both squarings and multiplications. The following section will describe the method for a reduced only two-power set.

2.3 Exponent Combination Method

Hereafter, we only consider the number of multiplications required to compute the r_i 's assuming that it is possible to calculate the squares only once by the previously described parallel method.

Denoting

$$a = \sum_i a_i 2^i, \quad b = \sum_i b_i 2^i$$

and assuming that $g^a \bmod p$ and $g^b \bmod p$ are to be computed, let $c = \sum_i c_i 2^i$.

If a and b are randomly chosen, one should expect that :

$$\sum_i a_i \equiv \sum_i b_i \equiv \frac{n}{2} \quad \text{and} \quad \sum_i c_i \equiv \frac{n}{4}.$$

Given the fact that :

$$w(a-c) + w(b-c) + w(c) \leq w(a) + w(b)$$

(where w denotes the Hamming weight), our strategy consists in computing :

$$\begin{cases} G_a = g^a \oplus c \bmod p \\ G_b = g^b \oplus c \bmod p \\ G_c = g^c \bmod p \end{cases}$$

to obtain

$$\begin{cases} r_a = G_a G_c \bmod p = g^a \bmod p \\ r_b = G_b G_c \bmod p = g^b \bmod p \end{cases}$$

The gain for a set of N signatures is therefore statistically $N(n/4 - 1)$ multiplications, which tends to 25% of the total multiplications required to generate a set of N signatures with the parallel square-&-multiply, if we simply apply this strategy to all signatures grouped by pairs as illustrated in the following table :

| operations | PSM | Exponent Combination |
|------------|------------------|------------------------|
| r_a | $\equiv n/2 - 1$ | $\equiv n/2 - n/4 - 1$ |
| r_b | $\equiv n/2 - 1$ | $\equiv n/2 - n/4 - 1$ |
| g^c | none | $\equiv n/4 - 1$ |
| Total | $\equiv n - 2$ | $\equiv 3n/4 - 3$ |

Table 1 : PSM and batch exponentiation performance

The computational effort required to generate the set

$$G = \{g^{k_i} \bmod p \mid i = 1, \dots, N\}$$

is about $N/2(3n/4 - 3) + n/N$ but implies to use $3N/2$ size(p)-bit registers which may not be practical in some situations. In the following we will present an optimization of our batch strategy and achieve comparison with Brickell, Gordon, McCurley and Wilson algorithm in [1], exhibiting when our strategy is more convenient and suitable.

3. Improvement and Performance

3.1 BGCW algorithm

The precomputation technique proposed by Brickell and al. in [1] is based on the following observation : in [11] it was proposed to precompute the set of g^{2^i} to reduce the computational effort increasing the storage amount. There is no reason to consider powers of 2.

The idea is to find a decomposition

$$k = \sum_{i=0}^{m-1} a_i x_i,$$

where $0 \leq a_i \leq h$ for $0 \leq i \leq m$, then we can compute

$$g^k = \prod_{d=1}^h c_d^d,$$

where $c_d = \prod_{a_i=d} g^{x_i}$.

The algorithm directly derived from this achieves the computation of g^k with only $m+h-2$ multiplications, but required that the values of $g^{x_i} \bmod p$ have been previously stored. To give an overview of the performances of this algorithm, one can remark that to generate a g^k with a 160-bit k , say for a Schnorr or DSS signature, using a base 16 representation for the numbers, the BGCW algorithm produces g^k at the cost of only 50 multiplications but requires also the storage of 40 n -bit numbers. Considering a base 32 notation, one can save a few storage (about 2 Kbytes rather than 2.5 previously) but the number of multiplications grows to 60. An improvement based on the notion of a *basic digit set* improves drastically the speed performance since the number of multiplications falls to about 36 but implies the storage of more than 200 numbers. The main drawback of the method is clearly the minimum storage capacity needed to achieve an exponentiation.

3.2 Optimizing the exponent combination

The total number of computations required to produce the set of g^k 's will depend mainly on the multiplications to be calculated since the squarings are done once for all.

Considering that we want to group α n -bit exponents from a large set of N values together rather than only joining them by pairs, the

main issue is to find the best combination strategy to reduce the multiplications to be done.

The number of multiplications per g^k is in average

$$M(n, \alpha) = \frac{n(2^\alpha - 1)}{\alpha 2^\alpha} + 2^\alpha - 1 - 1$$

and the squaring effort can be divided between all the k s.

The analysis of the function representation (see Annex 1) for usual exponent lengths (160-bit and 512-bit) give integer solutions which minimize this quantity. The number of registers required to achieve the computation increasing with the number of elements in the set (since to compute squarings once we must calculate all the g^k s together), the best values appear to be 4 for 160-bit values and 5 for 512-bit values. The final choice for a dedicated computation relies mainly on the memory capacity of the machine which generates the g^k s.

3.3 Performance Overview

The various strategies achieving several time and memory trade-offs are to be considered. The Tables 1 and 2 give the different trade-offs for 160-bit (Schnorr, DSS) and 512-bit (El-Gamal, Brickell-McCurley) exponents, considering that p is a 512-bit prime modulus.

| Set Size | Subset Size | Storage (Bytes) | #Multiplications per computation |
|----------|-------------|-----------------|----------------------------------|
| 2 | 2 | 316 | 141 |
| 3 | 3 | 652 | 103 |
| 4 | 2 | 568 | 101 |
| 4 | 4 | 1,324 | 84.5 |
| 12 | 3 | 2,416 | 63 |
| 12 | 4 | 3,844 | 57.83 |
| 36 | 3 | 7,120 | 54.11 |
| 36 | 4 | 11,404 | 48.94 |
| 108 | 4 | 34,084 | 45.98 |

Table 2 : performances with 160-bit exponent

| Set Size | Subset Size | Storage (Bytes) | #Multiplications per computation |
|----------|-------------|-----------------|----------------------------------|
| 2 | 2 | 448 | 449 |
| 3 | 3 | 960 | 323 |
| 4 | 4 | 1,984 | 255 |
| 5 | 5 | 4,032 | 216.6 |
| 60 | 3 | 17,984 | 160.87 |
| 60 | 4 | 28,864 | 135.53 |
| 60 | 5 | 47,680 | 122.73 |
| 200 | 5 | 158,784 | 116.76 |

Table 3 : performances with 512-bit exponent

Compare to the *Square-&-Multiply* and *BGCW* algorithms, the *Batch Exponentiation* strategy provides nice improvements to save computation and memory. The pairs grouping provides a 42 % gain on the total computation at the cost of only 2 n -bit registers, while the *BGCW* implies at least a 2 Kbytes storage. The technique is perfectly suited to low-memory environment where memory cost is high; furthermore, the exponent re-combination is easy to perform on any machine. On the other hand, *BGCW* algorithm is very efficient when at least a few Kbytes of permanent memory are available.

An adaptation of *Batch Exponentiation* tailored for high-speed transmission (see Annex 2) even provides greater improvements by sub-contracting the squaring effort to an external device assuming nothing on his security.

4. Conclusion, Extensions and Open Questions

We presented a strategy which can accelerate the generation of DLP-based signatures. The main characteristics of the batch-exponentiation technique presented in this article are summarized in the following table.

| scheme \Rightarrow effort | Batch (memory) | Batch (time) |
|--------------------------------|-------------------|-----------------|
| Schnorr | 141 N | 45.98 N |
| El-Gamal | 449 N | 116.76 N |

Table 4 : exponentiation performance

Several open questions appear interesting to explore to further improve the proposed strategies :

- For a power $k \in \mathbb{Z}_q$, try to find a such that $k' = a \cdot q + k$ where the hamming weight of k' is significantly small. Since computations are done modulo q , this transformation of k does not have any impact on the result itself but may well reduce the computation workload.
- Find an algorithm such that the construction of the subsets is optimal, that is the ordering of the k s results in as few computations as possible.

Acknowledgments

The authors would like to thank Jacques Stern for his pertinent remarks and gentle support. Furthermore, during the CRYPTO'95 Rump Session Yokua Tsuruoka pointed out that he described a similar method in [12] at JWS'93. We didn't know anything about this paper so that we can honestly consider Tsuruoka's paper as an independent prior discovery of the same algorithm.

References

- [1] E. Brickell, D. Gordon and K. McCurley, *Fast exponentiation with precomputation*, technical report no. SAND91-1836C, Sandia National Laboratories, Albuquerque, New-Mexico, October 1991.
- [2] A. Fiat, *Batch RSA*, Advances in cryptology: Proceedings of Crypto'89, LNCS, Springer-Verlag, 435, pp. 175-185.
- [3] FIPS PUB 186, February 1, 1993, *Digital Signature Standard*.
- [4] T. El-Gamal, *A public-key cryptosystem and a signature scheme based on discrete logarithms*, IEEE TIT, vol. IT-31:4, pp 469-472, 1985.
- [5] D. Naccache, D. M'raïhi, S. Vaudenay and D. Raphaël, *Can DSA be improved ? - Complexity Trade-Offs with the Digital Signature Standard*, Advances in cryptology: Proceedings of Eurocrypt'94, Perugia, LNCS 950, pp. 77-85, Springer-Verlag, 1995.

[6] J.-J. Quisquater and M. de Soete, *Speeding up smart-card RSA computation with Insecure Coprocessors*, Proceedings of Smart Cards 2000, 1989, pp. 191-197.

[7] P.J.N de Rooij, *On The Security of the Schnorr Scheme using Preprocessing*, Advances in cryptology: Proceedings of Eurocrypt'91, Brighton, LNCS 547, pp. 71-80, Springer-Verlag, 1991.

[8] P.J.N de Rooij, *On Schnorr's Preprocessing for Digital Signature Schemes*, Advances in cryptology: Proceedings of Eurocrypt'93, Lofthus, LNCS 765, pp. 435-439, Springer-Verlag, 1994.

[9] Ryo-Fuji-Hara, *Cipher Algorithms and Computational Complexity*, Bit 17 (1985), 954-959.

[10] C. Schnorr, *Efficient Identification and Signatures for Smart-Cards*, Advances in cryptology: Proceedings of Eurocrypt'89, Berlin, LNCS 435, pp. 239-252, Springer-Verlag, 1990.

[11] J. Stern and S.Vaudenay, *Personal Communication*, 1994.

[12] Y. Tsuruoka, *A Fast Algorithm on Addition Sequence*, JW-ISC'93, 1993

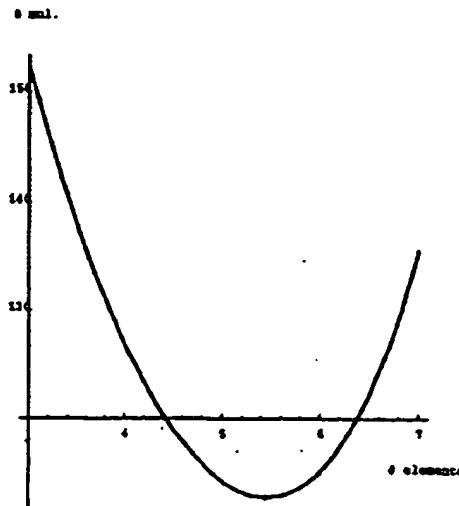


Figure 2 : 512-bit exponent - Best integer solution = 5

Appendix 2 : Sub-contracting squarings with high-speed transmission

Using a high-speed transmission interface such that of a PCMCIA card, one can subcontract the square computation and rely on the device with the same level of security.

The strategy remains the same, grouping the exponents and computing the g^k 's, except that the squaring are computed by a genuine device, assuming nothing on his tamper-resistance. The device that shall compute g^k values has a certificate C on the set of $\{g^{2^i} \bmod p \mid i \leq \text{size}(p)\}$ such as an iterative hashing of the whole set.

Denoting *Sender* the computing device in charge of the squaring effort and *Receiver* the exponentiation machine, the protocol is the following :

| Sender | Receiver |
|----------------------|---------------------------------|
| $s = g$ | $l = \emptyset$ |
| For $i = 0$ to $n-1$ | |
| Send s | Receive s (use if needed) |
| | $s = s^2 \bmod p$ |
| | $l = \text{SHA}(l \parallel s)$ |
| | If $C = l$ accept |

Appendix 1 : $M(n, \alpha) = \frac{n(2^\alpha - 1)}{\alpha 2^\alpha} + 2^{\alpha-1} - 1$

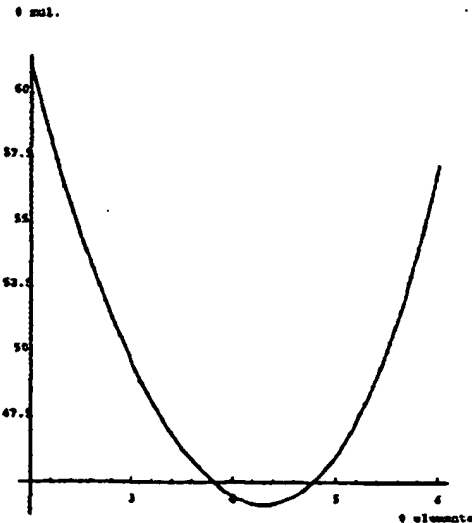


Figure 1 : 160-bit exponent - Best integer solution = 4

PicoDBMS: Scaling down database techniques for the smartcard

Philippe Pucheral¹, Luc Bouganim¹, Patrick Valduriez², Christophe Bobineau¹

¹ University of Versailles, PRISM Laboratory, Versailles, France;

E-mail: {philippe.pucheral;luc.bouganin;christophe.bobineau}@prism.uvsq.fr

² University Paris 6, LIP6 Laboratory, Paris, France; E-mail: patrick.valduriez@lip6.fr

Edited by A. El Abbadi, G. Schlageter, K.-Y. Whang. Received: 15 October 2000 / Accepted: 15 April 2001

Published online: 23 July 2001 – © Springer-Verlag 2001

Abstract. Smartcards are the most secure portable computing device today. They have been used successfully in applications involving money, and proprietary and personal data (such as banking, healthcare, insurance, etc.). As smartcards get more powerful (with 32-bit CPU and more than 1 MB of stable memory in the next versions) and become multi-application, the need for database management arises. However, smartcards have severe hardware limitations (very slow write, very little RAM, constrained stable memory, no autonomy, etc.) which make traditional database technology irrelevant. The major problem is scaling down database techniques so they perform well under these limitations. In this paper, we give an in-depth analysis of this problem and propose a PicoDBMS solution based on highly compact data structures, query execution without RAM, and specific techniques for atomicity and durability. We show the effectiveness of our techniques through performance evaluation.

Key words: Smartcard applications – PicoDBMS – Storage model – Execution model – Query optimization – Atomicity – Durability

1 Introduction

Smartcards are the most secure portable computing device today. The first smartcard was developed by Bull for the French banking system in the 1980s to significantly reduce the losses associated with magnetic stripe credit card fraud. Since then, smartcards have been used successfully around the world in various applications involving money, proprietary data, and personal data (such as banking, pay-TV or GSM subscriber identification, loyalty, healthcare, insurance, etc.). While today's smartcards handle a single issuer-dependent application, the trend is toward multi-application smartcards¹. Standards for multi-application support, like the JavaCard [36] and Microsoft's SmartCard for Windows [26], ensure that the card be universally accepted and be able to interact with several

service providers. This should make smartcards one of the world's highest-volume markets for semiconductors [14].

As smartcards become more and more versatile, multi-application, and powerful (32-bit processor, more than 1 MB of stable storage), the need for database techniques arises. Let us consider a health card storing a complete medical folder including the holder's doctors, blood type, allergies, prescriptions, etc. The volume of data can be important and the queries fairly complex (select, join, aggregate). Sophisticated access rights management using views and aggregate functions are required to preserve the holder's data privacy. Transaction atomicity and durability are also needed to enforce data consistency. More generally, database management helps to separate data management code from application code, thereby simplifying and making application code smaller. Finally, new applications can be envisioned, like computing statistics on a large number of cards, in an asynchronous and distributed way. Supporting database management on the card itself rather than on an external device is the only way to achieve very high security, high availability (anywhere, anytime, on any terminal), and acceptable performance.

However, smartcards have severe hardware limitations which stem from the obvious constraints of small size (to fit on a flexible plastic card and to increase hardware security) and low cost (to be sold in large volumes). Today's microcontrollers contain a CPU, memory – including about 96 kB of ROM, 4 kB of RAM, and up to 128 kB of stable storage like EEPROM – and security modules [39]. EEPROM is used to store persistent information; it has very fast read time (60–100 ns) comparable to old-fashion RAM but very slow write time (more than 1 ms/word). Following Moore's law for processor and memory capacities, smartcards will get rapidly more powerful. Existing prototypes, like Gemplus's Pinocchio card [16], bypass the current memory bottleneck by connecting an additional chip of 2 MB of Flash memory to the microcontroller. Although a significant improvement over today's cards, this is still very restricted compared to other portable, less secure, devices such as Personal Digital Assistants (PDA). Furthermore, smartcards are not autonomous, i.e., have no independent power supply, thereby precluding asynchronous and disconnected processing.

¹ Everyone would probably enjoy carrying far fewer cards.

These limitations (tiny RAM, little stable storage, very costly write, and lack of autonomy) make traditional database techniques irrelevant. Typically, traditional DBMS exploit significant amounts of RAM and use caching and asynchronous I/Os to reduce disk access overhead as much as possible. With the extreme constraints of the smartcard, the major problem is scaling down database techniques. While there has been much excellent work on scaling up to deal with very large databases, e.g., using parallelism, scaling down has not received much attention by the database research community. However, scaling down in general is becoming very important for commodity computing and is quite difficult [18].

Some DBMS designs have addressed the problem of scaling down. Light versions of popular DBMS like Sybase Adaptive Server Anywhere [37], Oracle 8i Lite [30] or DB2 Everywhere [20] have been primarily designed for portable computers and PDA. They have a small footprint which they obtain by simplifying and componentizing the DBMS code. However, they use relatively high RAM and stable memory and do not address the more severe limitations of smartcards. ISOL's SQLJava Machine DBMS [13] is the first attempt towards a smartcard DBMS while SCQL [24], the standard for smartcard database language, emerges. While both designs are limited to single select, they exemplify the strong interest for dedicated smartcard DBMS.

In this paper, we address the problem of scaling down database techniques and propose the design of what we call a PicoDBMS. This work is done in the context of a new project with Bull Smart Cards and Terminals. The design has been made with smartcard applications in mind, but its scope extends as well to any ultra-light computer device based on a secured monolithic chip. This paper makes the following contributions:

- We analyze the requirements for a PicoDBMS based on a typical healthcare application and justify its minimal functionality.
- We give an in-depth analysis of the problem by considering the smartcard hardware trends and derive design principles for a PicoDBMS.
- We propose a new pointer-based storage model that integrates data and indices in a unique compact data structure.
- We propose query execution techniques which handle complex query plans (including joins and aggregates) with no RAM consumption.
- We propose transaction techniques for atomicity and durability that reduce the logging cost to its lowest bound and enable a smartcard to participate in distributed transactions.
- We show the effectiveness of each technique through performance evaluation.

This paper is an extended version of [7]. In particular, the section on transaction management is new. The paper is organized as follows. Section 2 illustrates the use of take-away databases in various classes of smartcard applications and presents in more detail the requirements of the health card application. Section 3 analyzes the smartcard hardware constraints and gives the problem definition. Sections 4–6 present and assess the PicoDBMS' storage model, query execution model, and transaction model, respectively. Section 7 concludes.

2 Smartcard applications

In this section, we discuss the major classes of emerging smartcard applications and their database requirements. Then, we illustrate these requirements in further detail with the health card application, which we will use as reference example in the rest of the paper.

2.1 Database management requirements

Table 1 summarizes the database management requirements of the following typical classes of smartcard applications:

- *Money and identification*: examples of such applications are credit cards, e-purse, SIM for GSM, phone cards, transportation cards. They are representative of today's applications, with very few data (typically the holder's identifier and some status information). Querying is not a concern and access rights are irrelevant since cards are protected by PIN-codes. Their unique database management requirement is update atomicity.
- *Downloadable databases*: these are predefined packages of confidential data (e.g., diplomatic, military or business information) that can be downloaded on the card – for example, before traveling – and be accessed from any terminal. Data availability and security are the major concerns here. The volume of data can be important and the queries complex. The data are typically read-only.
- *User environment*: the objective is to store in a smartcard an extended profile of the card's holder including, among others, data regarding the computing environment (PC's configuration, passwords, cookies, bookmarks, software licenses, etc.), an address book as well as an agenda. The user environment can thus be dynamically recovered from the profile on any terminal. Queries remain simple, as data are not related. However, some of the data are highly private and must be protected by sophisticated access rights (e.g., the card's holder may want to share a subset of her/his address book or bookmark list with a subset of persons). Transaction atomicity and durability are also required.
- *Personal folders*: personal folders may be of a different nature: scholastic, healthcare, car maintenance history, loyalty. They roughly share the same requirements, which we illustrate next with the healthcare example. Note that queries involving data issued from different folders can make sense. For instance, one may be interested in discovering associations between some disease and the scholastic level of the card holder. This raises the interesting issue of maintaining statistics on a population of cards or mining their content asynchronously.

2.2 The health card application

The health card is very representative of personal folder applications and has strong database requirements. Several countries (France, Germany, USA, Russia, Korea, etc.) are developing healthcare applications on smartcards [11]. The initial idea was to give to each citizen a smartcard containing her/his identification and insurance data. As smartcard storage capacity increases, the information stored in the card can be

Table 1. Typical applications' profiles

| Applications | Volume | Select/project | Join | Group by / Distinct | Access rights / views | Atomicity | Durability | Statistics |
|------------------------|--------|----------------|------|---------------------|-----------------------|-----------|------------|------------|
| Money & identification | tiny | | | | | | ✓ | |
| Downloadable DB | high | ✓ | ✓ | ✓ | | | | |
| User environment | medium | ✓ | | | ✓ | ✓ | ✓ | |
| Personal folder | high | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

extended to the holder's doctors, emergency data (blood type, allergies, vaccination, etc.), surgical operations, prescriptions, insurance data and even links to heavier data (e.g., X-ray examination, scanner images, etc.) stored on hospital servers. Different users may query, modify, and create data in the holder's folder: the doctors who consult the patient's past records and prescribe drugs, the surgeons who perform exams and operations, the pharmacists who deliver drugs, the insurance agents who refund the patient, public organizations which maintain statistics or study the impact of drugs correlation in population samples, and finally the holder her/himself.

We can easily observe that: (i) the amount of data is significant (more in terms of cardinality than in terms of volume because most data can be encoded); (ii) queries can be rather complex (e.g., a doctor asks for the last antibiotics prescribed to the patient); (iii) sophisticated access rights management using views and aggregate functions are highly required (e.g., a statistical organization may access aggregate values only but not the raw data); (iv) atomicity must be preserved (e.g., when the pharmacist delivers drugs); and (v) durability is mandatory, without compromising data privacy (logged data stored outside the card must be protected).

One may wonder whether the holder's health data ought to be stored in a smartcard or in a centralized database. The benefit of distributing the healthcare database on smartcards is threefold. First, health data must be made highly available (anywhere, anytime, on any terminal, and without requiring a network connection). Second, storing sensitive data on a centralized server may damage privacy. Third, maintaining a centralized database is fairly complex due to the variety of data sources. Assuming the health data is stored in the smartcard, the next question is why the aforementioned database capabilities need to be hosted in the smartcard rather than the terminals. The answer is again availability (the data must be exploited on any terminal) and privacy. Regarding privacy, since the data must be confined in the chip, so must the query engine and the view manager. As the smartcard is the unique trusted part of the system, access rights and transaction management cannot be delegated to an untrusted terminal.

3 Problem formulation

In this section, we make clear the smartcard constraints in order to derive design rules for the PicoDBMS and state the problem. Our analysis is based on the characteristics of both

existing smartcard products and current prototypes [16, 39], and thus, should be valid for a while. We also discuss how the main constraints of the smartcard will evolve in a near future.

3.1 Smartcard constraints

Current smartcards include in a monolithic chip, a 32 bits RISC processor at about 30 MIPS, memory modules (of about 96 kB of ROM, 4 kB of static RAM, and 128 kB of EEPROM), security components (to prevent tampering), and take their electrical energy from the terminal [39]. ROM is used to store the operating system, the JavaCard virtual machine, fixed data, and standard routines. RAM is used as working memory for maintaining an execution stack and calculating results. EEPROM is used to store persistent information. EEPROM has very fast read time (60–100 ns/word) comparable to old-fashion RAM, but a dramatically slow write time (more than 1 ms/word).

The main constraints of current smartcards are therefore: (i) the very limited storage capacity; (ii) the very slow write time in EEPROM; (iii) the extremely reduced size of the RAM; (iv) the lack of autonomy; and (v) a high security level that must be preserved in all situations. These constraints strongly distinguish smartcards from any other computing devices, including lightweight computers like PDA.

Let us now consider how hardware advances can impact on these constraints, in particular, memory size. Current smartcards rely on a well-established and slightly out-of-date hardware technology (0.35 μm) in order to minimize the production cost (less than five dollars) and increase security [34]. Furthermore, up to now, there was no real need for large memories in smartcard applications such as the holder's identification. According to major smartcard providers, the market pressure generated by emerging large storage demanding applications will lead to a rapid increase of the smartcard storage capacity. This evolution is however constrained by the smartcard tiny die size fixed to 25 mm² in the ISO standard [23], which pushes for more integration. This limited size is due to security considerations (to minimize the risk of physical attack [5]) and practical constraints (e.g., the chip should not break when the smartcard is flexed). Another solution to relax the storage limit is to extend the smartcard storage capacity with external memory modules. This is being done by Gemplus which recently announced Pinocchio [16], a smartcard equipped with 2 MB of Flash memory linked to the microcontroller by a bus. Since hardware security can no longer be provided on this memory, its content must be either non-sensitive or encrypted.

Another important issue is the performance of stable memory. Possible alternatives to the EEPROM are Flash memory and Ferroelectric RAM (FeRAM) [15] (see Table 2 for performance comparisons). Flash is more compact than EEPROM and represents a good candidate for high capacity smartcards [16]. However, Flash banks need to be erased before writing, which is extremely slow. This makes Flash memory appropriate for applications with a high read/write ratio (e.g., address books). FeRAM is undoubtedly an interesting option for smartcards as read and write times are both fast. Although its theoretical foundation was set in the early 1950s, FeRAM is just emerging as an industrial solution. Therefore, FeRAM is expensive, less secure than EEPROM or Flash, and its integration with traditional technologies (such as CPUs) remains an

Table 2. Performance of stable memories for the smartcard

| Memory type | EEPROM | FLASH | FeRAM |
|----------------------|------------------------------|------------------------------|---|
| Read time (/word) | 60 to 150 ns | 70 to 200 ns | 150 to 200 ns |
| Write time (/word) | 1 to 5 ms | 5 to 10 μ s | 150 to 200 ns |
| Erase time (/bank) | None | 500 to 800 ms | None |
| Lifetime (*) (/cell) | 10 ⁵ write cycles | 10 ⁵ erase cycles | 10 ¹⁰ to 10 ¹² write cycles |

* A memory cell can be overwritten a finite number of time.

issue. Thus FeRAM could be considered a serious alternative only in the very long term [15].

Given these considerations, we assume in this paper a smartcard with a reasonable stable storage area (a few megabytes of EEPROM²) and a small RAM area (some kilobytes). Indeed, there is no clear interest in having a large RAM area, given that the smartcard is not autonomous, thus precluding asynchronous write operations. Moreover, more RAM means less EEPROM as the chip size is limited.

3.2 Impact on the PicoDBMS architecture

We now analyze the impact of the smartcard constraints on the PicoDBMS architecture, thus justifying why traditional database techniques, and even lightweight DBMS techniques, are irrelevant. The smartcard's properties and their impact are:

- **Highly secure:** smartcard's hardware security makes it the ideal storage support for private data. The PicoDBMS must contribute to the data security by providing access right management and a view mechanism that allows complex view definitions (i.e., supporting data composition and aggregation). The PicoDBMS code must not present security holes due to the use of sophisticated algorithms³.
- **Highly portable:** the smartcard is undoubtedly the most portable personal computer (the wallet computer). The data located on the smartcard are thus highly available. They are also highly vulnerable since the smartcard can be lost, stolen or accidentally destroyed. The main consequence is that durability cannot be enforced locally.
- **Limited storage resources:** despite the foreseen increase in storage capacity, the smartcard will remain the lightest representative of personal computers for a long time. This means that specific storage models and execution techniques must be devised to minimize the volume of persistent data (i.e., the database) and the memory consumption during execution. In addition, the functionalities of the PicoDBMS must be carefully selected and their implementation must be as light as possible. The lightest the PicoDBMS, the biggest the onboard database.
- **Stable storage is main memory:** smartcard stable memory provides the read speed and direct access granularity of a main memory. Thus, a PicoDBMS can be considered as a *main memory DBMS (MMDBMS)*. However the dramatic cost of writes distinguishes a PicoDBMS from a traditional MMDBMS. This impacts on the storage and access

² Considering Flash instead of EEPROM will not change our conclusions. It will just exacerbate them.

³ Most security holes are the results of software bugs [34].

methods of the PicoDBMS as well as the way transaction atomicity is achieved.

- **Non-autonomous:** compared to other computers, the smartcard has no independent power supply, thereby precluding disconnected and asynchronous processing. Thus, all transactions must be completed while the card is inserted in a terminal (unlike PDA, write operations cannot be cached in RAM and reported on stable storage asynchronously).

3.3 Problem statement

To summarize, our goal is to design a PicoDBMS including the following components:

- **Storage manager:** manages the storage of the database and the associated indices.
- **Query manager:** processes query plans composed of select, project, join, and aggregates.
- **Transaction manager:** enforces the ACID properties and participates to distributed transactions.
- **Access right manager:** provides access rights on base data and on complex user-defined views.

Thus, the PicoDBMS hosted in the chip provides the minimal subset of functionality that is strictly needed to manage in a secure way the data shared by all onboard applications. Other components (e.g., the GUI, a sort operator, etc.) can be hosted in the terminal or be dynamically downloaded when needed, without threatening security. In the rest of this paper, we concentrate on the components which require non-traditional techniques (storage manager, query manager, and transaction manager) and ignore the access right manager for which traditional techniques can be used.

When designing the PicoDBMS's components, we must follow several design rules derived from the smartcard's properties:

- **Compactness rule:** minimize the size of data structures and the PicoDBMS code to cope with the limited stable memory area (a few megabytes).
- **RAM rule:** minimize the RAM usage given its extremely limited size (some kilobytes).
- **Write rule:** minimize write operations given their dramatic cost (≈ 1 ms/word).
- **Read rule:** take advantage of the fast read operations (≈ 100 ns/word).
- **Access rule:** take advantage of the low granularity and direct access capability of the stable memory for both read and write operations.
- **Security rule:** never externalize private data from the chip and minimize the algorithms' complexity to avoid security holes.

4 PicoDBMS storage model

In this section, following the design rules for a PicoDBMS, we discuss the storage issues and propose a very compact model based on a combination of flat storage, domain storage, and ring storage. We also evaluate the storage cost of our storage model.

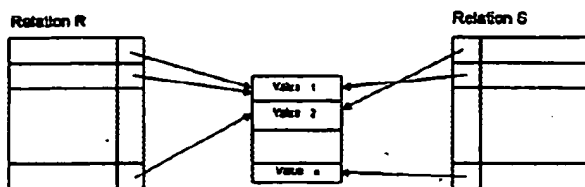


Fig. 1. Domain storage

4.1 Flat storage

The simplest way to organize data is *flat storage (FS)*, where tuples are stored sequentially and attribute values are embedded in the tuples. Although it does not impose it, the SCQL standard [24] considers FS as the reference storage model for smartcards. The main advantage of FS is access locality. However, in our context, FS has two main drawbacks:

- *Space consuming*: while normalization rules preclude attributes conjunction redundancy to occur, they do not avoid attribute value duplicates (e.g., the attribute *Doctor.Specialty* may contain many duplicates).
- *Inefficient*: in the absence of index structures, all operations are computed sequentially. While this is convenient for old fashion cards (some kilobytes of storage and a mono-relation select operator), this is no longer acceptable for future cards where storage capacity is likely to exceed 1 MB and queries can be rather complex.

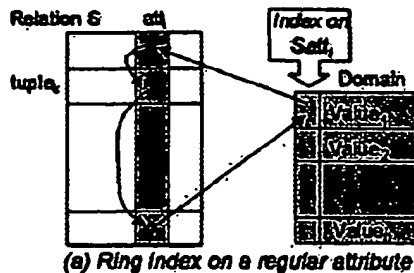
Adding index structures to FS may solve the second problem while worsening the first one. Thus, FS alone is not appropriate for a PicoDBMS.

4.2 Domain storage

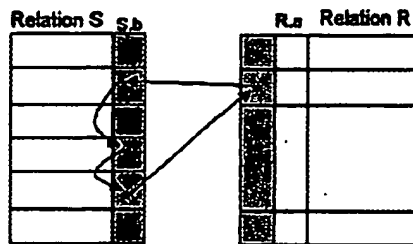
Based on the critique of FS, it follows that a PicoDBMS storage model should guarantee both data and index compactness. Let us first deal with data compactness. Since locality is no longer an issue in our context, pointer-based storage models inspired by MMDBMS [3, 27, 31] can help reducing the data storage cost. The basic idea is to preclude any duplicate value from occurring. This can be achieved by grouping values in domains (sets of unique values). We call this model *domain storage (DS)*. As shown in Fig. 1, tuples reference their attribute values by means of pointers. Furthermore, a domain can be shared among several attributes. This is particularly efficient for enumerated types, which vary on a small and determined set of values⁴.

One may wonder about the cost of tuple creation, update, and deletion since they may generate insertion and deletion of values in domains. While these actions are more complex than their FS counterpart, their implementation remains more efficient in the smartcard context, simply because the amount of data to be written is much smaller. To amortize the slight overhead of domain storage, we only store by domain all large attributes (i.e., greater than a pointer size) containing duplicates. Obviously, attributes with no duplicates (e.g., keys) need

⁴ Compression techniques can be advantageously used in conjunction with DS to increase compactness [17].



(a) Ring Index on a regular attribute



(b) Ring Index on a foreign key attribute

Fig. 2. Ring storage

not be stored by domain but with FS. Variable-size attributes – generally larger than a pointer – can also be advantageously stored in domains even if they do not contain duplicates. The benefit is not storage savings but memory management simplicity (all tuples of all relations become fixed-size) and log compactness (see Sect. 6).

4.3 Ring storage

We now address index compactness along with data compactness. Unlike disk-based DBMS that favor indices which preserve access locality, smartcards should make intensive use of secondary (i.e., pointer-based) indices. The issue here is to make these indices as compact as possible. Let us first consider select indices. A select index is typically made of two parts: a collection of values and a collection of pointers linking each value to all tuples sharing it. Assuming the indexed attribute varies on a domain, the index's collection of values can be saved since it exactly corresponds to the domain extension. The extra cost incurred by the index is then reduced to the pointers linking index values to tuples.

Let us go one step further and get these pointers almost for free. The idea is to store these *value-to-tuple* pointers in place of the *tuple-to-value* pointers within the tuples (i.e., pointers stored in the tuples to reference their attribute values in the domains). This yields to an index structure which makes a ring from the domain values to the tuples. Hence, we call it *ring index* (see Fig. 2a). However, the ring index can also be used to access the domain values from the tuples and thus serve as data storage model. Thus we call *ring storage (RS)* the storage of a domain-based attribute indexed by a ring. The index storage cost is reduced to its lowest bound, that is, one pointer per domain value, whatever the cardinality of the indexed relation. This important storage saving is obtained at the price of extra work for projecting a tuple to the corresponding attribute since retrieving the value of a ring stored attribute means traversing

on average half of the ring (i.e., up to reaching the domain value).

Join indices [40] can be treated in a similar way. A join predicate of the form $(R.a = S.b)$ assumes that $R.a$ and $S.b$ vary on the same domain. Storing both $R.a$ and $S.b$ by means of rings leads to defining a join index. In this way, each domain value is linked by two separate rings to all tuples from R and S sharing the same join attribute value. However, most joins are performed on key attributes, $R.a$ being a primary key and $S.b$ being the foreign key referencing $R.a$. In our model, key attributes are not stored by domain but with FS. Nevertheless, since $R.a$ is the primary key of R , its extension forms precisely a domain, even if not stored outside of R . Since attributes $S.b$ take their values in $R.a$'s domain, they reference $R.a$ values by means of pointers. Thus, the domain-based storage model naturally implements for free a *unidirectional join index* from $S.b$ to $R.a$ (i.e., each S tuple is linked by a pointer to each R tuple matching with it). If traversals from $R.a$ to $S.b$ need to be optimized too, a *bi-directional join index* is required. This can be simply achieved by defining a ring index on $S.b$. Figure 2b shows the resulting situation where each R tuple is linked by a ring to all S tuples matching with it and vice versa. The cost of a bi-directional join index is restricted to a single pointer per R tuple, whatever the cardinality of S . Note that this situation resembles the well-known Codasyl model.

4.4 Storage cost evaluation

Our storage model combines FS, DS, and RS. Thus, the issue is to determine the best storage for each attribute. If the attributes need not be indexed, the choice is obviously between FS and DS. Otherwise, the choice is between RS and FS with a traditional index. Thus, we compare the storage cost for a single attribute, indexed or not, for each alternative. We introduce the following parameters:

- *CardRel*: cardinality of the relation holding the attribute.
- *a*: average length of the attribute (expressed in bytes).
- *p*: pointer size (3 bytes will be required to address "large" memory of future cards).
- *S*: selectivity factor of the attribute. $S = CardDom/CardRel$, where *CardDom* is the cardinality of the attribute domain extension (in all models). *S* measures the redundancy of the attribute (i.e., the same attribute value appears in $1/S$ tuples).

$Cost(FS) = CardRel * a$ // attribute storage cost in
// the relation

$Cost(DS) = CardRel * p$ // attribute storage cost in
// the relation
+ $S * CardRel * a$ // values storage cost in
// the domain

$Cost(Indexed.FS) = Cost(FS)$ // flat attribute storage cost
+ $S * CardRel * a$ // value storage cost in the
// index
+ $CardRel * p$ // pointer storage cost in
// the index

$Cost(RS) = Cost(DS)$ // domain-based attribute
// storage cost
+ $S * CardRel * p$ // pointer storage cost in
// the index

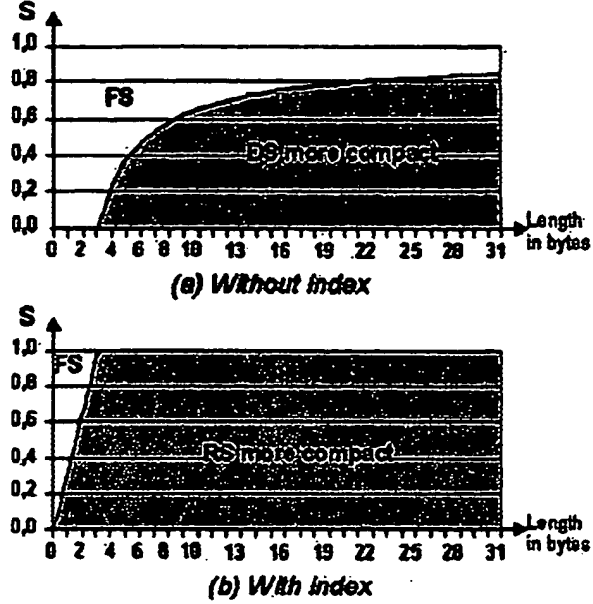


Fig. 3. Storage models' tradeoff

The cost equality between FS and DS gives: $S = (a-p)/a$.
The cost equality between Indexed.FS and RS gives:

$$S = a/p$$

Figure 3a shows the different values of *S* and *a* for which FS and DS are equivalent. Thus, each curve divides the plan into a gain area for FS (above the curve) and a gain area for DS (under the curve). For values of *a* less than 3 (i.e., the size of a pointer), FS is obviously always more compact than DS. For higher values of *a*, DS becomes rapidly more compact than FS except for high values of *S*. For instance, considering $S = 0.5$, that is the same value is shared by only two tuples, DS outperforms FS for all *a* larger than 6 bytes. The higher *a* and the lower *S*, the better DS. The benefit of DS is thus particularly important for enumerated type attributes. Figure 3b compares Indexed.FS with RS. The superiority of RS is obvious, except for 1- and 2-byte-long key attributes. Thus, Figs. 3a and 3b are guidelines for the database designer to decide how to store each attribute, by considering its size and selectivity.

5 Query processing

Traditional query processing strives to exploit large main memory for storing temporary data structures (e.g., hash tables) and intermediate results. When main memory is not large enough to hold some data, state-of-the-art algorithms (e.g., hybrid hash join [33]) resort to materialization on disk to avoid memory overflow. These algorithms cannot be used for a PicoDBMS because:

- Given the write rule and the lifetime of stable memory, writes in stable memory are proscribed, even for temporary materialization.

- Dedicating a specific RAM area does not help since we cannot estimate its size a priori. Making it small increases the risk of memory overflow, thereby leading to writes in stable memory. Making it large reduces the stable memory area, already limited in a smartcard (RAM rule). Moreover, even a large RAM area cannot guarantee that query execution will not produce memory overflow [9].
- State-of-the-art algorithms are quite sophisticated, which precludes their implementation in a PicoDBMS whose code must be simple, compact, and secure (compactness and security rules).

To solve this problem, we propose query processing techniques that do not use any working RAM area nor incur any writes in stable memory. In the following, we describe these techniques for simple and complex queries, including aggregation and remove duplicates. We show the effectiveness of our solution through a performance analysis.

5.1 Basic query execution without RAM

We consider the execution of *SPJ* (*Select/Project/Join*) queries. Query processing is classically done in two steps. The query optimizer first generates an "optimal" *query execution plan* (*QEP*). The QEP is then executed by the query engine which implements an *execution model* and uses a library of relational operators [17]. The optimizer can consider different shapes of QEP: *left-deep*, *right-deep* or *bushy trees* (see Fig. 4). In a left-deep tree, operators are executed sequentially and each intermediate result is materialized. On the contrary, right-deep trees execute operators in a pipeline fashion, thus avoiding intermediate result materialization. However, they require materializing in memory all left relations. Bushy trees offer opportunities to deal with the size of intermediate results and memory consumption [38].

In a PicoDBMS, the query optimizer should not consider any of these execution trees as they incur materialization. The solution is to only use pipelining with *extreme right-deep trees* where all the operators (including select) are pipelined. As left operands are always base relations, they are already materialized in stable memory, thus allowing us to execute a plan with no RAM consumption. Pipeline execution can be easily achieved using the well-known *Iterator Model* [17]. In this model, each operator is an *iterator* that supports three procedure calls: *open* to prepare an operator for producing an item, *next* to produce an item, and *close* to perform final clean-up. A *QEP* is activated starting at the root of the operator tree and progressing towards the leaves. The dataflow in the model is demand-driven: a child operator passes a tuple to its parent node in response to a *next* call from the parent.

Let us now detail how select, project, and join are performed. These operators can be executed either sequentially or with a ring index. Given the access rule, the use of indices seems always to be the right choice. However, extreme right-deep trees allow us to speed-up a single select on the first base relation (e.g., *Drug.type* in our example), but using a ring index on the other selected attributes (e.g., *Visit.date*) may slow down execution as the rings need to be traversed to retrieve their value. Project operators are pushed up to the tree since no materialization occurs. Note that the final project incurs

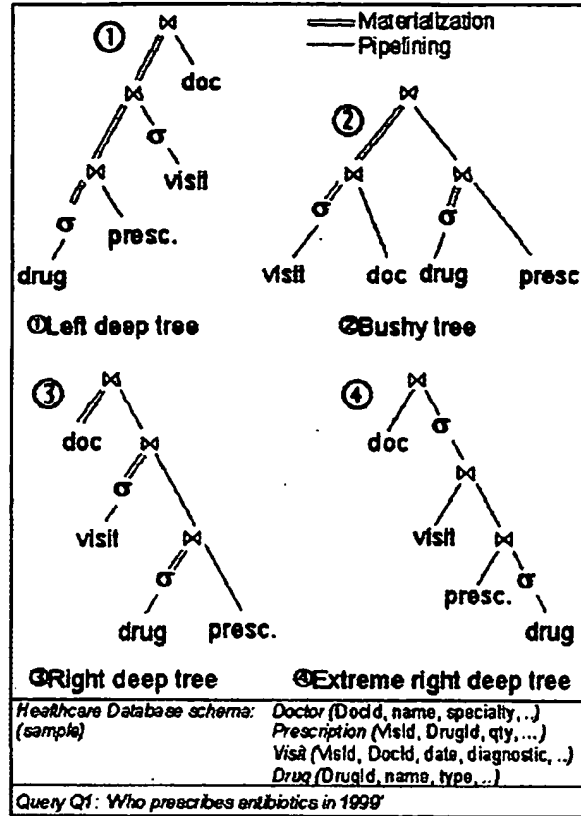


Fig. 4. Several execution trees for query Q1

an additional cost in case of ring attributes. Without indices, joining relations is done by a nested-loop algorithm since no other join technique can be applied without ad hoc structures (e.g., hash tables) and/or working area (e.g., sorting). The cost of indexed joins depends on the way indices are traversed. Consider the indexed join between *Doctor* (n tuples) and *Visit* (m tuples) on their key attribute. Assuming a unidirectional index, the join cost is proportional to $n * m$ starting with *Doctor* and to m starting with *Visit*. Assuming now a bi-directional index, the join cost becomes proportional to $n + m$ starting with *Doctor* and to $m^2/2n$ starting with *Visit* (retrieving the doctor associated to each visit incurs traversing half of a ring in average). In the latter case, a naïve nested loop join can be more efficient if the ring cardinality is greater than the target relation cardinality (i.e., when $m > n^2$). In that case, the database designer must clearly choose a unidirectional index between the two relations.

5.2 Complex query execution without RAM

We now consider the execution of aggregate, sort, and duplicate removal operators. At first glance, pipeline execution is not compatible with these operators which are classically performed on materialized intermediate results. Such materialization cannot occur either in the smartcard due to the RAM rule or in the terminal due to the security rule. Note that sorting can be done in the terminal since the output order of the

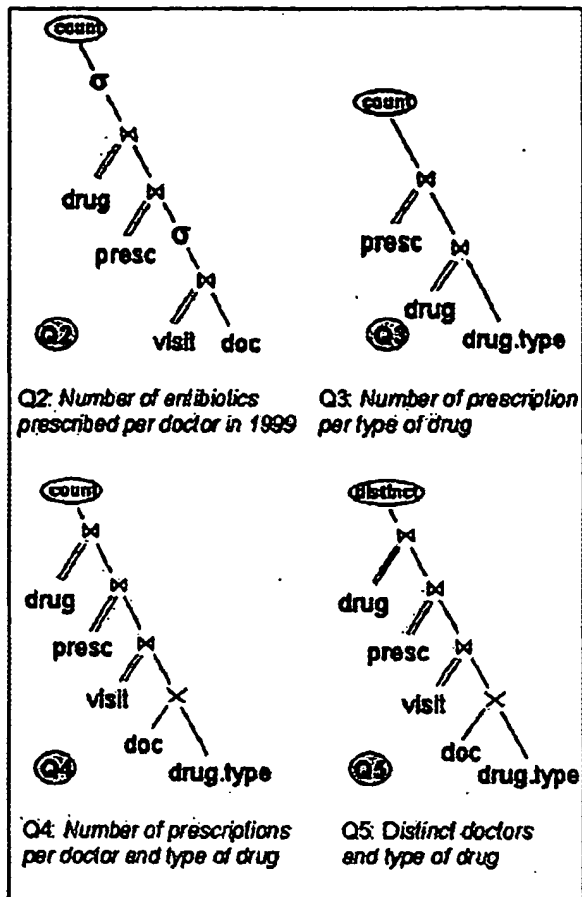


Fig. 5. Four 'complex' query execution plans

result tuples is not significant, i.e., depends on the DBMS algorithms.

We propose a solution to the above problem by exploiting two properties: (i) aggregate and duplicate removal can be done in pipeline if the incoming tuples are still grouped by distinct values; and (ii) pipeline operators are order-preserving since they consume (and produce) tuples in the arrival order. Thus, enforcing an adequate consumption order at the leaf of the execution tree allows pipelined aggregation and duplicate removal. For instance, the extreme right-deep tree of Fig. 4 delivers the tuples naturally grouped by *Drug.id*, thus allowing group queries on that attribute.

Let us now consider query Q2 of Fig. 5. As pictured, executing Q2 in pipeline requires rearranging the execution tree so that relation *Doctor* is explored first. Since *Doctor* contains distinct doctors, the tuples arriving to the *count* operator are naturally grouped by doctors.

The case of Q3 is harder. As the data must be grouped by *type of drugs* rather than by *Drug.id*, an additional join is required between relation *Drug* and domain *drug.type*. Domain values being unique, this join produces the tuples in the adequate order. If domain *Drug.type* does not exist, an operator must be introduced to sort relation *Drug* in pipeline. This can be done by performing n passes on *Drug* where n is the number of distinct values of *Drug.type*.

The case of Q4 is even trickier. The result must be grouped on two attributes (*Doctor.id* and *Drug.type*), introducing the need to start the tree with both relations! The solution is to insert a Cartesian product operator at the leaf of the tree in order to produce tuples ordered by *Doctor.id* and *Drug.type*. In this particular case, the query response time should be approximately n times greater than the same query without the 'group by' clause, where n is the number of distinct *types of drugs*.

Q5 retrieves the distinct couples of *doctor* and *type of prescribed drugs*. This query can be made similar to Q4 by expressing the distinct clause as an aggregate without function (i.e., the query "select distinct a_1, \dots, a_n from ..." is equivalent to "select a_1, \dots, a_n from ... group by a_1, \dots, a_n "). The unique difference is that the computation for a given group, i.e., (distinct result tuple) can stop as soon as one tuple has been produced.

5.3 Query optimization

Heuristic optimization is attractive. However, well-known heuristics such as processing select and project first do not work here. Using extreme right-deep trees makes the former impractical and invalidates the latter. Heuristics for join ordering are even more risky considering our data structures. Conversely, there are many arguments for an exhaustive search of the best plan. First, the search space is limited since: (i) there is a single algorithm for each operator, depending on the existing indices; (ii) only extreme right-deep trees are considered; and (iii) typical queries will not involve many relations. Second, exhaustive search using depth-first algorithms do not consume any RAM. Finally, exhaustive algorithms are simple and compact (even if they iterate a lot). Under the assumption that query optimization is required in a PicoDBMS, the remarks above strongly argue in favor of an exhaustive search strategy.

5.4 Performance evaluation

Our proposed query engine can handle fairly complex queries, taking advantage of the read and access rules⁵ while satisfying the compactness, write, RAM, and security rules. We now evaluate whether the PicoDBMS performance matches the smartcard application's requirements, that is, any query issued by the application can be performed in reasonable time (i.e., may not exceed the user's patience). Since the PicoDBMS code's simplicity is an important consideration to conform to the compactness and security rules, we must also evaluate which acceleration techniques (i.e., ring indices, query optimization) are really mandatory. For instance, an accelerator reducing the response time from 10ms to 1ms is useless in the smartcard context⁶. Thus, unlike traditional performance evaluation, our major concern is on absolute rather than relative performance.

⁵ With traditional DBMS, such techniques will induce so many disk accesses that the system would thrash!

⁶ With traditional DBMS, such acceleration can improve the transactional throughput.

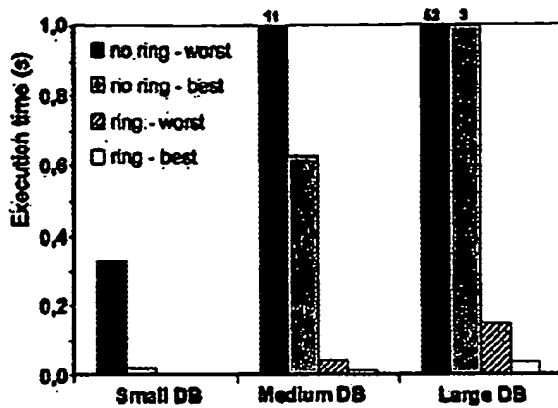


Fig. 6. Performance results for Q1

Evaluating absolute response time is complex in the smartcard environment because all platform parameters (e.g., processor speed, caching strategy, RAM, and EEPROM speed) strongly impact on the measurements⁷. Measuring the performance of our PicoDBMS on Bull's smartcard technology is attractive but introduces two problems. First, Bull's smartcards compatible with database applications are still prototypes [39]. Second, we are interested in providing the most general conclusions (i.e., as independent as possible of smartcard architectures). Therefore, we prefer to measure our query engine on two old-fashioned computers (a PC 486/25 Mhz and a Sun SparcStation 1+) which we felt roughly similar to forthcoming smartcard architectures. For each computer, we vary the system parameters (clock frequency, cache) and perform the experimentation tests. The performance ratios between all configurations were roughly constant (i.e., whatever the query), the slowest configuration (Intel 486 with no cache) performing eight times worse than the fastest (RISC with cache). In the following, we present response times for the slowest architecture to check the viability of our solutions in the worst environment.

We generated three instances of a simplified healthcare database: the *small*, *medium*, and *large* databases containing, respectively, (10, 30, 50) doctors, (100, 500, 1,000) visits, (300, 2,000, 5,000) prescriptions, and (40, 120, 200) drugs. Although we tested several queries, we describe below only the two most significant. Query Q1, which contains three joins and two selects on *Visit* and *Drug* (with selectivities of 20% and 5%), is representative of medium-complexity queries. Query Q4, which performs an aggregate on two attributes and requires the introduction of a Cartesian product, is representative of complex queries. For each query, we measure the performance for all possible query execution plans, excluding those which induce additional Cartesian product, varying the storage choices (with and without select and join ring indices). Figures 6 and 7 show the results for both best and worst plans on databases built with or without join indices.

Considering SPJ queries, the PicoDBMS performance clearly matches the application's requirements as soon as join rings are used. Indeed, the performance with join rings is at

⁷ With traditional DBMS, very slow disk access allows us to ignore finer parameters.

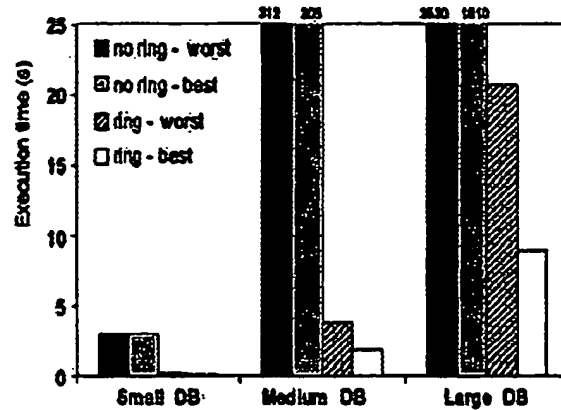


Fig. 7. Performance results for Q4

most 146 ms for the largest database and with the worst execution plan. With small databases, all the acceleration techniques can be discarded, while with larger ones, join rings remain necessary to obtain good response time. In that case, the absolute gain (110 ms) between the best and the worst plan does not justify the use of a query optimizer.

The performance of aggregate queries is clearly the worst because they introduce a Cartesian product at the leaf of the execution tree. Join rings are useful for medium and large databases. With large databases, the optimizer turns out to be necessary since the worst execution plan with join rings achieves a rather long response time (20.6 s).

The influence of ring indices for selects (not shown) is insignificant. Depending on the selectivity, it can bring slight improvement or overhead on the results. Although it may achieve an important relative speed-up for the select itself, the absolute gain is not significant considering the small influence of select on the global query execution cost (which is not the case in disk-based DBMS). Select ring indices are, however, useful for queries with aggregates or duplicate removal, that can result in a join between a relation and the domain attribute. In that case, the select index plays the role of a join index, thereby generating a significant gain on large relations and large domains.

Thus, this performance evaluation shows that our approach is feasible and that join indices are mandatory in all cases while query optimization turns out to be useful only with large databases and complex queries.

6 Transaction management

Like any data server, a PicoDBMS must enforce the well-known transactional ACID properties [8] to guarantee the consistency of the local data it manages as well as be able to participate in distributed transactions. We discuss below these properties with respect to a PicoDBMS.

- *Atomicity: local atomicity* means that the set of actions performed by the PicoDBMS on a transaction's behalf is made persistent following the *all or nothing* scheme. *Global atomicity*: this means that all data servers – including the PicoDBMS – accessed by a distributed transaction

agree on the same transaction outcome (either commit or rollback). The distinguishing features of a PicoDBMS regarding atomicity are no demarcation between main memory and persistent storage, the dramatic cost of writes, and the fact that they cannot be deferred.

- **Consistency:** this property ensures that the actions performed by the PicoDBMS satisfy all integrity constraints defined on the local data. Considering that traditional integrity constraint management can be used, we do not discuss it any further.
- **Isolation:** this property guarantees the serializability of concurrent executions. A PicoDBMS manages personal data and is typically single-user⁸. Furthermore, smartcard operating systems do not even support multithreading. Therefore, isolation is useless here.
- **Durability:** durability means that committed updates are never lost whatever the situation (i.e., even in case of a media failure). Durability cannot be enforced locally by the PicoDBMS because the smartcard is more likely to be stolen, lost or destroyed than a traditional computer. Indeed, mobility and smallness play against safety. Consequently, durability must be enforced through the network. The major issue is then preserving the privacy of data while delegating the durability to an external agent.

The remainder of this section addresses local atomicity, global atomicity, and durability.

6.1 Local atomicity

There are basically two ways to perform updates in a DBMS. The updates are either performed on *shadow objects* that are atomically integrated in the database at commit time or done *in place* (i.e., the transaction updates the shared copy of the database objects) [8]. We discuss these two traditional models below.

- **Shadow update:** This model is rarely employed in disk-based DBMSs because it destroys data locality on disk and increases concurrent updates on the catalog. In a PicoDBMS, disk locality and concurrency are not a concern. This model has been shown to be convenient for smartcards equipped with a small Flash memory [25]. However, it is poorly adapted to pointer-based storage models like RS since the object location changes at every update. In addition, the cost incurred by shadowing grows with the memory size. Indeed, either the granularity of the shadow objects increases or the paths to be duplicated in the catalog become longer. In both cases, the writing cost – which is the dominant factor – increases.
- **Update in-place:** write-ahead logging (WAL) [8] is required in this model to undo the effects of an aborted transaction. Unfortunately, the relative cost of WAL is much higher in a PicoDBMS than in a traditional disk-based DBMS which uses buffering to minimize I/Os. In a smartcard, the log must be written for each update since each update becomes immediately persistent. This roughly doubles the cost of writing.

⁸ Even if the data managed by the PicoDBMS are shared among multiple users (e.g., as in the healthcare application), the PicoDBMS serves a single user at a time.

Despite its drawbacks, *update in-place* is better suited than *shadow update* for a PicoDBMS because it accommodates pointer-based storage models and its cost is insensitive to the rapid growth of stable memory capacity. We also propose two optimizations to *update in-place*:

- **Pointer-based logging:** traditional WAL logs the values of all modified data. RS allows a finer granularity by logging pointers in place of values. The smallest the log records, the cheapest the WAL. The logging process must consider two types of information:
- **Values:** in case of a tuple update, the log record must contain the tuple address and the old attribute values, that is a pointer for all RS stored attributes and a regular value for FS stored attributes. In case of a tuple insertion or deletion, assuming each tuple header contains a status bit (i.e., dead or alive), only the tuple address has to be logged in order to recover its state.
- **Rings:** tuple insertion, deletion, and update (of a ring attribute) modify the structure of each ring traversing the corresponding tuple t . Since a ring is a circular chain of pointers, recovering its state means recovering the *next* pointer of t 's predecessor (let us call it t_{pred}). The information to restore in $t_{pred}.next$ is either t 's address if t has been updated or deleted, or $t.next$ if t has been inserted. t 's address already belongs to the log (see above) and $t.next$ does not have to be logged since t 's content still exists in stable storage at recovery time. The issue is how to identify t_{pred} at recovery time. Logging this information can be saved at the price of traversing the whole ring starting from t , until reaching t again. Thus, ring recovery comes for free in terms of logging.
- **Garbage-collecting values:** insertion and deletion of domain values (domain values are never modified) should be logged as any other updates. This overhead can be avoided by implementing a deferred garbage collector that destroys all domain values no longer referenced by any tuple. Garbage-collecting a domain amounts to execute an ad hoc semi-join operator between the domain and all relations varying on it which discards the domain values that do not match⁹. The benefit of this solution is threefold: (i) the lazy deletion of unreferenced values does not entail the storage model coherency; (ii) garbage-collecting domain values is required anyway by RS (even in the absence of transaction control); and (iii) a deferred garbage-collector can be implemented without reference counters, thereby saving storage space. The deferred garbage collector cannot work in the background since smartcards do not yet support multi-threading. The most pragmatic solution is to launch it manually when the card is nearly full. An alternative to this manual procedure is to execute the garbage collector automatically at each card connection on a very small subset of the database (so that its cost remains hidden to the user). Garbage-collecting the database in such an incremental way is straightforward since domain values are examined one after the other.

⁹ Unlike reachability algorithms that start from the persistent roots and need marking [6], the proposed garbage-collector starts from the persistent leaves (i.e., the domain values) and exploits them one after the other, in a pipelined fashion (thus, it conforms to the RAM rule).

The update in-place model along with pointer-based logging and deferred garbage-collector reduces logging cost to its lowest bound, that is, a tuple address for inserted and deleted tuples, and the values of updated attributes (again, a pointer for DS and RS stored attributes).

6.2 Global atomicity

Global atomicity is traditionally enforced by an *atomic commitment protocol (ACP)*. The most well known and widely used ACP is 2PC [8]. While extensively studied [19] and standardized [21, 29, 41], 2PC suffers from the following weaknesses in our context:

- *Need for a standard prepared state*: any server must externalize the standard *Xa* interface [41] to participate to 2PC. Unfortunately, ISO defines a transactional interface for smartcards but it does not cover distributed transactions [24]. In addition, participating to 2PC requires building a local prepared state that consumes valuable resources.
- *Disconnection means aborting*: a smartcard can be extracted from its terminal or its mobile host (e.g., a cellular phone) can be temporarily unreachable during 2PC. A participant's disconnection leads 2PC to abort the transaction even if all its operations have been successfully executed.
- *Badly adapted to moving participants*: the 2PC incurs two message rounds to commit a transaction. Considering the high cost of wireless communication, the overhead is significant for mobile terminals equipped with a smartcard reader (e.g., PDA, cellular phones).

As its name indicates, 2PC has two phases: the *voting* phase and the *decision* phase. The voting phase is the means by which the coordinator checks whether or not the participants can locally guarantee the ACID properties of the distributed transaction. The decision is *commit* if all participants vote *yes* and *abort* otherwise. Thus, the voting phase introduces an uncertainty period at transaction termination that leads to the aforementioned drawbacks.

Variations of *one-phase commit* protocols (*1PC*) have been recently proposed [2, 4, 35]. As stated in [2], 1PC eliminates the voting phase of 2PC by enforcing the following properties on the participant's behavior: (1) all operations are acknowledged before the 1PC is launched; (2) there are no deferred integrity constraints; (3) all participants are ruled by a rigorous concurrency control scheduler; and (4) all updates are logged on stable storage before 1PC is launched. These assumptions guarantee, respectively, the A, C, I, D properties before the ACP is launched. Then, the ACP reduces to a single phase, that is broadcasting the coordinator's decision to all participants (this decision is *commit* if all transaction's operations have been successfully executed and *abort* otherwise). If a crash or a disconnection precludes a participant from conforming to this decision, the corresponding transaction branch is simply forward recovered (potentially at the next reconnection). While the assumptions on the participant's behavior seem constraining in the general case, they are quite acceptable in the smartcard context [10]. Property (1) is common to all ACPs and is enforced by the ISO7816 standard [22]; property (2) conforms to the fact that PicoDBMS have lighter capabilities

than full-fledged DBMS; and property (3) is satisfied by definition since smartcards do not support parallel executions. Property (4) is discussed in Sect. 6.3.

Eliminating the voting phase of the ACP solves altogether the three aforementioned problems. However, one may wonder about the interoperability between transaction managers and data managers supporting different protocols (either 1PC or 2PC). We have shown in [1] that the participation of legacy (i.e., 2PC compliant) data managers in 1PC is straightforward. Conversely, the participation of 1PC compliant data managers (e.g., a smartcard) in the 2PC can be achieved by associating a *log agent* to each participant. The role of the log agent is twofold. First, it manages the data manager's part of the 1PC's coordinator log, forces it to stable storage during the 2PC prepare phase, and exploits it if the transaction branch needs to be forward-recovered. Second, it translates the 2PC interface into that of 1PC. The log agent can be located on the terminal, so that the benefit of 1PC is lost for the terminal but it is preserved for the smartcard.

6.3 Durability

Most 1PC protocols assume that the coordinator is in charge of logging all participants' updates before triggering the ACP (all these protocols belong to the coordinator log family). *Coordinator log* [35] and *implicit yes vote* [4] assume that the participants piggyback their log records on the acknowledgment messages of each operation while *coordinator logical log* [2] assumes that the coordinator logs all operations sent to each participant. In all cases, the durability of the distributed transaction relies on the coordinator log. Thus, 1PC is a means by which global atomicity and durability can be solved altogether, at the same price.

Two issues remain to be solved: (i) where to store the coordinator log; and (ii) how to preserve the security rule, that is, how to make the log content as secure as the data stored in the smartcard. Since the log must sustain any kind of failure, it must be stored on the network by a trustee server (e.g., a public organism, a central bank, the card issuer, etc.). If some transactions are executed in disconnected mode (e.g., on a mobile terminal), the durability will be effective only at the time the terminal reconnects to the network. Protecting the log content against attacks imposes encryption. The way encryption is performed depends on the model of logging. If the coordinator log is fed by the log records piggybacked by the participants, the smartcard can encrypt them with an algorithm based on a private key (e.g., DES [28]). Otherwise (i.e., if the *coordinator logical log* scheme is selected), the smartcard can provide the coordinator with a public key that will be used by the coordinator itself to encrypt its log [32].

6.4 Transaction cost evaluation

The goal of this section is to approximate the time required by a representative update transaction. The objective is to confirm whether or not the write performance of smartcards assumed in this paper is acceptable for database applications like health cards. To this end, we estimate the time required to create a tuple in a relation, including the creation of domain values,

the insertion of the tuple in the rings potentially defined on this relation and the log time. Let us introduce the following parameters, in addition to those already defined in Sect. 4.4:

- *nbAttFS*: number of FS stored attributes
- *nbAttDS*: number of DS stored attributes
- *nbAttRS*: number of RS stored attributes
- *w*: size of a word (4 bytes in a 32-bit card)
- *t*: time to write one word in stable storage (5 ms in the worst case)

```
Cost(insertTuple) =
  ((nbAttFS*a + nbAttDS*p + nbAttRS*p)/w) // ①
  + (nbAttRS + nbAttDS) * S * [a/w] // ②
  + nbAttRS * [p/w] // ③
  + [p/w] // ④
  ) * t // ⑤
```

- ① Tuple size
- ② Domain values size. $S \approx$ probability to create a new domain value
- ③ Ring pointers to be updated
- ④ Log record size
- ⑤ Write time

Let us consider a representative transaction executed on the healthcare. This transaction inserts a new tuple in *Doctor* and *Visit* and five tuples in *Prescription* and *Drug*. This is somehow a worst case for this application in the sense that the visited doctor is a new one and prescribes five new drugs. The considered attribute distribution is as follows:

| | |
|---------------------|---|
| <i>Doctor</i> | (<i>nbAttFS</i> =3, <i>nbAttDS</i> =4, <i>nbAttRS</i> =0), |
| <i>Visit</i> | (<i>nbAttFS</i> =2, <i>nbAttDS</i> =3, <i>nbAttRS</i> =2), |
| <i>Prescription</i> | (<i>nbAttFS</i> =1, <i>nbAttDS</i> =1, <i>nbAttRS</i> =2), |
| <i>Drug</i> | (<i>nbAttFS</i> =2, <i>nbAttDS</i> =4, <i>nbAttRS</i> =0). |

The average attribute length *a* is fixed to 10 bytes. Figure 8 plots the update transaction execution time depending on *S* (*S* = 0 means that all attribute values already exist in the domains, while *S* = 1 means that all these values need be inserted in the domains).

The figure is self-explanatory. Note that the logging cost represents less than 3% of the total cost. This simple analysis shows that the time expected for this kind of transaction (less than 1 s) is clearly compatible with the healthcare application's requirements.

7 Conclusion

As smartcards become more and more versatile, multi-application, and powerful, the need for database techniques arises. However, smartcards have severe hardware limitations which make traditional database technology irrelevant. The major problem is scaling down database techniques so they perform well under these limitations. In this paper, we addressed this problem and proposed the design of a PicoDBMS, concentrating on the components which require non-traditional techniques (storage manager, query manager, and transaction manager).

This paper makes several contributions. First, we analyzed the requirements for a PicoDBMS based on a healthcare application which is representative of personal

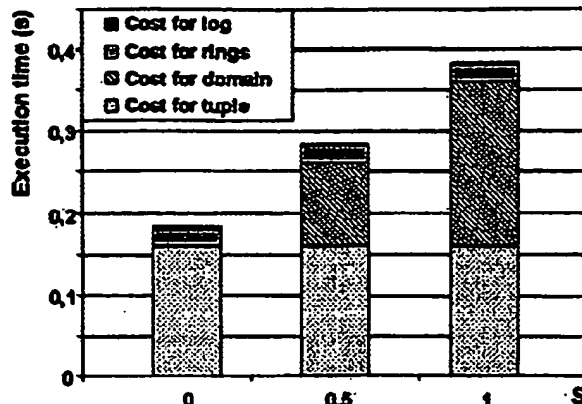


Fig. 8. Performance of a typical update transaction

folder applications and has strong database requirements. We showed that the minimal functionality should include select/project/join/aggregate, access right management, and views as well as transaction's atomicity and durability.

Second, we gave an in-depth analysis of the problem by considering the smartcard hardware trends. Based on this analysis, we assumed a smartcard with a reasonable stable memory of a few megabytes and a small RAM of some kilobytes, and we derived design rules for a PicoDBMS architecture.

Third, we proposed a new highly compact storage model that combines flat storage (FS), domain storage (DS), and ring storage (RS). Ring storage reduces the indexing cost to its lowest bound. Based on performance evaluation, we derived guidelines to decide the best way to store an attribute.

Fourth, we proposed query processing techniques which handle complex query plans with no RAM consumption. This is achieved by considering extreme right-deep trees which can pipeline all operators of the plan including aggregates. We also argued that, if query optimization is needed, the strategy should be exhaustive search. We measured the performance of our execution model with an implementation of our query engine on two old-fashioned computers which we configured to be similar to forthcoming smartcard architectures. We showed that the resulting performance matches the smartcard application's requirements.

Finally, we proposed techniques for transaction atomicity and durability. Local atomicity is achieved through update in-place with two optimizations which exploit the storage model: pointer-based logging and garbage collection of domain values. Global atomicity and durability are enforced by 1PC which is easily applicable in the smartcard context and more efficient than 2PC. We showed that the performance of typical update transactions is acceptable for representative applications like the health card.

This work is done in the context of a new project with Bull Smart Cards and Terminals. The next step is to port our PicoDBMS prototype on Bull's smartcard new technology, called *OverSoft* [12], and to assess its functionality and performance on real-world applications. To this end, a benchmark dedicated to PicoDBMS must be set up. We also plan to address open issues such as protected logging for durability, query execution on encrypted data (e.g., stored in an external Flash), and statistics maintenance on a population of cards.

References

1. Abdallah M., Bobineau C., Guerraoui R., Pucheral P.: Specification of the transaction service. Esprit project OpenDREAMS-II n° 25262, Deliverable n° R13, 1998
2. Abdallah M., Guerraoui R., Pucheral P.: One-phase commit: does it make sense? Int. Conf. on Parallel and Distributed Systems (ICPADS), 1998
3. Ammann A., Hanrahan M., Krishnamurthy R.: design of a memory resident DBMS. IEEE COMPCON, 1985
4. Al-Houmaili Y., Chrysanthos P.K.: Two-phase commit in gigabit-networked distributed databases. Int. Conf. on Parallel and Distributed Computing Systems (PDCS), 1995
5. Anderson R., Kuhn M.: Tamper resistance – a cautionary note. USENIX Workshop on Electronic Commerce, 1996
6. Amsaleg L., Franklin M.J., Gruber O.: Efficient incremental garbage collection for client-server object database systems. Int. Conf. on Very Large Databases (VLDB), 1995
7. Bobineau C., Bouganim L., Pucheral P., Valduriez P.: PicoDBMS: scaling down database techniques for the smartcard (Best Paper Award). Int. Conf. on Very Large Databases (VLDB), 2000
8. Bernstein P.A., Hadzilacos V., Goodman N.: Concurrency control and recovery in database systems. Addison-Wesley, Reading, Mass., USA, 1987
9. Bouganim L., Kapitkaia O., Valduriez P.: Memory-adaptive scheduling for large query execution. Int. Conf. on Information and Knowledge Management (CIKM), 1998
10. Bobineau C., Pucheral P., Abdallah M.: A unilateral commit protocol for mobile and disconnected computing. Int. Conf. on Parallel and Distributed Computing Systems (PDCS), 2000
11. van Bommel F.A., Sembritzki J., Buettner H.-G.: Overview on healthcard projects and standards. Health Cards Int. Conf., 1999
12. Bull S.A.: Bull unveils iSimplify! the personal portable portal. Available at: <http://www.bull.com:80/bull.news/>
13. Carrasco L.C.: RDBMS's for Java cards? What a senseless ideal. Available at: www.sqlmachine.com, 1999
14. DataQuest.: Chip card market and technology charge ahead. MSAM-WW-DP-9808, 1998
15. Dipert B.: FRAM: Ready to ditch niche? EDN Access Magazine, Cahners, London, 1997
16. Gemplus.: SIM Cards: From kilobytes to megabytes. Available at: www.gemplus.fr/about/pressroom/, 1999
17. Graefe G.: Query evaluation techniques for large databases. ACM Comput Surv, 25(2), 1993
18. Graefe G.: The new database imperatives. Int. Conf. on Data Engineering (ICDE), 1998
19. Gray J., Reuter A.: Transaction processing. Concepts and Techniques. Morgan Kaufmann, San Francisco, 1993
20. IBM Corporation.: DB2 Everywhere – administration and application programming guide. IBM Software Documentation, SC26-9675-00, 1999
21. International Standardization Organization (ISO): Information technology - open systems interconnection - distributed transaction processing. ISO/IEC 10026, 1992
22. International Standardization Organization (ISO): Integrated circuit(s) cards with contacts – part 3: electronic signal and transmission protocols. ISO/IEC 7816-3, 1997
23. International Standardization Organization (ISO): Integrated circuit(s) cards with contacts – part 1: physical characteristics. ISO/IEC 7816-1, 1998
24. International Standardization Organization (ISO): Integrated circuit(s) cards with contacts – part 7: interindustry commands for structured card query language (SCQL). ISO/IEC 7816-7, 1999
25. Lecomte S., Trane P.: Failure recovery using action log for smartcards transaction based system. IEEE Online Testing Workshop, 1997
26. Microsoft Corporation.: Windows for smartcards toolkit for visual basic 6.0. Available at: www.microsoft.com/windowsce/smartcard/, 2000
27. Missikoff M., Scholl M.: Relational queries in a domain based DBMS. ACM SIGMOD Int. Conf. on Management of Data, 1983
28. National Institute of Standards and Technology.: Announcing the Data Encryption Standard (DES). FIPS PUB 46-2, 1993
29. Object Management Group.: Object transaction service. Document 94.8.4, OMG editor, 1994
30. Oracle Corporation.: Oracle 8i Lite - Oracle Lite SQL reference. Oracle documentation, A73270-01, 1999
31. Pucheral P., Thévenin J.M., Valduriez P.: Efficient main memory data management using the DBGraph storage model. Int. Conf. on Very Large Databases (VLDB), 1990
32. RSA Laboratories.: PKCS # 1: RSA Encryption Standard. RSA Laboratories Technical Note, v.1.5, 1993
33. Schneider D., DeWitt D.: A performance evaluation of four parallel join algorithms in a shared-nothing multiprocessor environment. ACM-SIGMOD Int. Conf., 1989
34. Schneier B., Shostack A.: Breaking up is hard to do: modeling security threats for smart cards. USENIX Symposium on Smart Cards, 1999
35. Stamos J., Cristian F.: A low-cost atomic commit protocol. IEEE Symposium on Reliable Distributed Systems, 1990
36. Sun Microsystems.: JavaCard 2.1 application programming interface specification. JavaSoft documentation, 1999
37. Sybase Inc.: Sybase adaptive server anywhere reference. CT75KNA, 1999
38. Shekita E., Young H., Tan K.L.: Multi-join optimization for symmetric multiprocessors. Int. Conf. on Very Large Data Bases (VLDB), 1993
39. Tual J.-P.: MASSC: a generic architecture for multiapplication smart cards. IEEE Micro J, N° 0272-1739/99, 1999
40. Valduriez P.: Join indices. ACM Trans. Database Syst, 12(2), 1987
41. X/Open.: Distributed transaction processing: reference model. X/Open Guide, Version 3. G307., X/Open Company Limited, 1996

Implementations for Coalesced Hashing

Jeffrey Scott Vitter
Brown University

The coalesced hashing method is one of the faster searching methods known today. This paper is a practical study of coalesced hashing for use by those who intend to implement or further study the algorithm. Techniques are developed for tuning an important parameter that relates the sizes of the address region and the cellar in order to optimize the average running times of different implementations. A value for the parameter is reported that works well in most cases. Detailed graphs explain how the parameter can be tuned further to meet specific needs. The resulting tuned algorithm outperforms several well-known methods including standard coalesced hashing, separate (or direct) chaining, linear probing, and double hashing. A variety of related methods are also analyzed including deletion algorithms, a new and improved insertion strategy called varied-insertion, and applications to external searching on secondary storage devices.

CR Categories and Subject Descriptors: D.2.8 [Software Engineering]: Metrics—*performance measures*; E.2 [Data]: Data Storage Representations—*hash-table representations*; F.2.2 [Analysis of Algorithms and Problem Complexity]: Nonnumerical Algorithms and Problems—*sorting and searching*; H.2.2 [Database Management]: Physical Design—*access methods*; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*search process*

General Terms: Algorithms, Design, Performance, Theory

Additional Key Words and Phrases: analysis of algorithms, coalesced hashing, hashing, data structures, databases, deletion, asymptotic analysis, average-case, optimization, secondary storage, assembly language

This research was supported in part by a National Science Foundation fellowship and by National Science Foundation grants MCS-77-23738 and MCS-81-05324.

Author's Present Address: Jeffrey Scott Vitter, Department of Computer Science, Box 1910, Brown University, Providence, RI 02912.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission. © 1982 ACM 0001-0782/82/1200-0911 \$00.75.

1. Introduction

One of the primary uses today for computer technology is information storage and retrieval. Typical searching applications include dictionaries, telephone listings, medical databases, symbol tables for compilers, and storing a company's business records. Each package of information is stored in computer memory as a *record*. We assume there is a special field in each record, called the *key*, that uniquely identifies it. The job of a searching algorithm is to take an input K and return the record (if any) that has K as its key.

Hashing is a widely used searching technique because no matter how many records are stored, the average search times remain *bounded*. The common element of all hashing algorithms is a predefined and quickly computed *hash function*

$$\text{hash: (all possible keys)} \rightarrow \{1, 2, \dots, M\}$$

that assigns each record to a *hash address* in a uniform manner. (The problem of designing hash functions that justify this assumption, even when the distribution of the keys is highly biased, is well-studied [7, 2].) Hashing methods differ from one another by how they resolve a *collision* when the hash address of the record to be inserted is already occupied.

This paper investigates the *coalesced hashing* algorithm, which was first published 22 years ago and is still one of the faster known searching methods [16, 7]. The total number of available storage locations is assumed to be *fixed*. It is also convenient to assume that these locations are contiguous in memory. For the purpose of notation, we shall number the hash table slots $1, 2, \dots, M'$. The first M slots, which serve as the range of the hash function, constitute the *address region*. The remaining $M' - M$ slots are devoted solely to storing records that collide when inserted; they are called the *cellar*. Once the cellar becomes full, subsequent colliders must be stored in empty slots in the address region and, thus, may trigger more collisions with records inserted later.

For this reason, the search performance of the coalesced hashing algorithm is very sensitive to the relative sizes of the address region and cellar. In Sec. 4, we apply the analytic results derived in [10, 11, 13] in order to optimize the ratio of their sizes, $\beta = M/M'$, which we call the *address factor*. The optimizations are based on two performance measures: the number of probes per search and the running time of assembly language versions. There is no unique best choice for β —the optimum address factor depends on the type of search, the number of inserted records, and the performance measure chosen—but we shall see that the compromise choice $\beta \approx 0.86$ works well in many situations. The method can be further turned to meet specific needs.

Section 5 shows that this tuned method dominates several popular hashing algorithms including standard coalesced hashing (in which $\beta = 1$), separate (or direct)

chaining, linear probing, and double hashing. The last three sections deal with variations and different implementations for coalesced hashing including deletion algorithms, alternative insertion methods, and external searching on secondary storage devices.

This paper is designed to provide a comprehensive treatment of the many practical issues concerned with the implementation of the coalesced hashing method. Readers interested in the theoretical justification of the results in this paper can consult [10, 11, 13, 14, 1].

2. The Coalesced Hashing Algorithm

The algorithm works like this: Given a record with key K , the algorithm searches for it in the hash table, starting at location $hash(K)$ and following the links in the chain. If the record is present in the table, then it is found and the search is *successful*; otherwise, the end of the chain is reached and the search is *unsuccessful*. For simplicity, we assume that the record is inserted whenever the search ends unsuccessfully, according to the following rule: If position $hash(K)$ is empty, then the record is stored at that location; else, it is placed in the largest-numbered empty slot in the table and is linked to the end of the chain. This has the effect of putting the first $M' - M$ colliders into the cellar.

Coalesced hashing is a generalization of the well-known separate (or direct) chaining method. The separate chaining method halts with overflow when there is no more room in the cellar to store a collider. The example in Fig. 1(a) can be considered to be an example

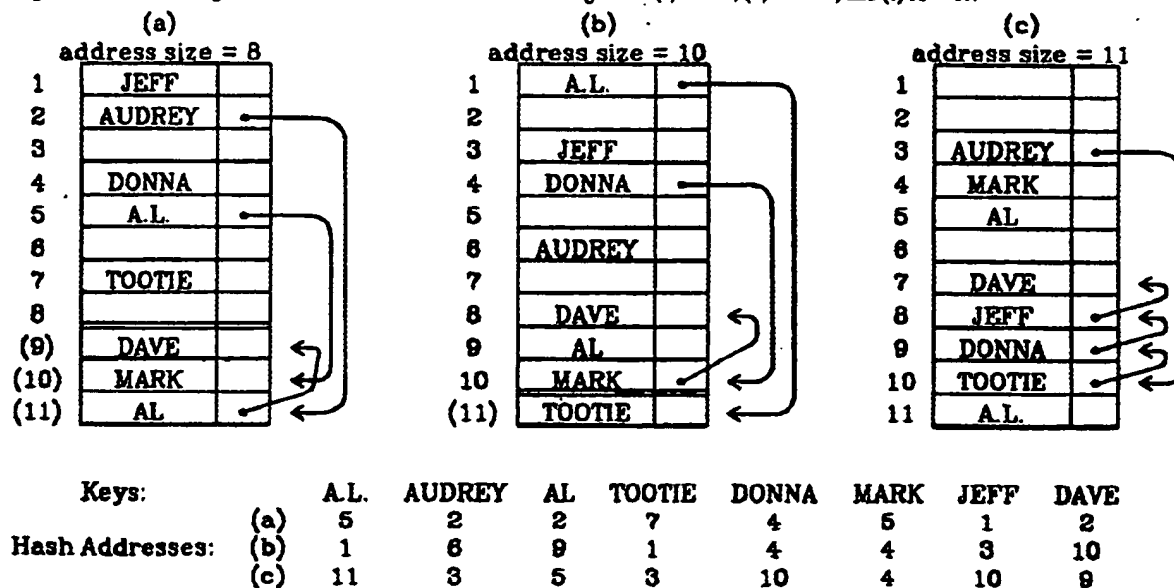
of both coalesced hashing and separate chaining, because the cellar is large enough to store the three colliders.

Figures 1(b) and 1(c) show how the two methods differ. The cellar contains only one slot in the example in Fig. 1(b). When the key MARK collides with DONNA at slot 4, the cellar is already full. Separate chaining would report overflow at this point. The coalesced hashing method, however, stores the key MARK in the largest-numbered empty space (which is location 10 in the address region). This causes a later collision when DAVE hashes to position 10, so DAVE is placed in slot 8 at the end of the chain containing DONNA and MARK. The method derives its name from this "coalescing" of records with different hash addresses into single chains.

The average number of probes per search shows marked improvement in Fig. 1(b), even though coalescing has occurred. Intuitively, the larger address region spreads out the records more evenly and causes fewer collisions, i.e., the hash function can be thought of as "shooting" at a bigger target. The cellar is now too small to store these fewer colliders, so it overflows. Fortunately, this overflow occurs late in the game, and the pileup phenomenon of coalescing is not significant enough to counteract the benefits of a larger address region. However, in the extreme case when $M = M' = 11$ and there is no cellar (which we call *standard coalesced hashing*), coalescing begins too early and search time worsens (as typified by Figure 1(c)). Determining the optimum address factor $\beta = M/M'$ is a major focus of this paper.

The first order of business before we can start a detailed study of the coalesced hashing method is to formalize the algorithm and to define reasonable measures of search performance. Let us assume that each

Fig. 1. Coalesced hashing, $M' = 11$, $N = 8$. The sizes of the address region are (a) $M = 8$, (b) $M = 10$, and (c) $M = 11$.



average # probes per successful search: (a) $12/8 = 1.5$, (b) $11/8 = 1.375$, (c) $14/8 = 1.75$.

of the M' contiguous slots in the coalesced hash table has the following organization:

| | | | |
|---------------------------------|-------|--------------|--------|
| E M P T Y | KEY | other fields | $LINK$ |
|---------------------------------|-------|--------------|--------|

For each value of i between 1 and M' , $EMPTY[i]$ is a one-bit field that denotes whether the i th slot is unused, $KEY[i]$ stores the key (if any), and $LINK[i]$ is either the index to the next spot in the chain or else the null value 0.

The algorithms in this article are written in the English-like style used by Knuth in order to make them readily understandable to all and to facilitate comparisons with the algorithms contained in [7, 4, 12]. Block-structured languages, like PL/I and Pascal, are good for expressing complicated program modules; however, they are not used here, because hashing algorithms are so short that there is no reason to discriminate against those who are not comfortable with such languages.

Algorithm C (Coalesced hashing search and insertion). This algorithm searches an M' -slot hash table, looking for a given key K . If the search is unsuccessful and the table is not full, then K is inserted.

The size of the address region is M ; the hash function *hash* returns a value between 1 and M (inclusive). For convenience, we make use of slot 0, which is always empty. The global variable R is used to find an empty space whenever a collision must be stored in the table. Initially, the table is empty, and we have $R = M' + 1$; when an empty space is requested, R is decremented until one is found. We assume that the following initializations have been made before any searches or insertions are performed: $M \leftarrow \lceil \beta M' \rceil$, for some constant $0 < \beta \leq 1$; $EMPTY[i] \leftarrow \text{true}$, for all $0 \leq i \leq M'$; and $R \leftarrow M' + 1$.

- C1. [Hash.] Set $i \leftarrow \text{hash}(K)$. (Now $1 \leq i \leq M$.)
- C2. [Is there a chain?] If $EMPTY[i]$, then go to step C6. (Otherwise, the i th slot is occupied, so we will look at the chain of records that starts there.)
- C3. [Compare.] If $K = KEY[i]$, the algorithm terminates successfully.
- C4. [Advance to next record.] If $LINK[i] \neq 0$, then set $i \leftarrow LINK[i]$ and go back to step C3.
- C5. [Find empty slot.] (The search for K in the chain was unsuccessful, so we will try to find an empty table slot to store K .) Decrease R one or more times until $EMPTY[R]$ becomes true. If $R = 0$, then there are no more empty slots, and the algorithm terminates with overflow. Otherwise, append the R th cell to the chain by setting $LINK[i] \leftarrow R$; then set $i \leftarrow R$.
- C6. [Insert new record.] Set $EMPTY[i] \leftarrow \text{false}$, $KEY[i] \leftarrow K$, $LINK[i] \leftarrow 0$, and initialize the other fields in the record. ■

In this paper, we concern ourselves with measuring the *searching phase* of Algorithm C and ignore for the most part the insertion time in steps C5 and C6. (The time for step C5 is not significant, because the total number of times R is decremented over the course of all the insertions cannot be more than the number of inserted records; hence, the amortized expected number of decrements is at most 1. The decrementing operation can also be done in parallel with steps C1–C4.) Our primary measure of search performance is the *number of probes per search*, which is the number of different table slots that are accessed while searching. In Algorithm C, this quantity is equal to

$$\max\{1, \text{number of times step C3 is performed}\}$$

For example, in Fig. 1(b), the unsuccessful searches for keys A.L. and TOOTIE (immediately prior to their insertions) each took one probe, while a successful search for DAVE would take two probes.

The average performance of the algorithm is obtained by assuming that all searches and insertions are random. The Appendix contains a discussion of the probability model as well as the formulas for the expected number of probes in unsuccessful and successful searches.

3. Assembly Language Implementation

Even though probe-counting gives us a good idea of search performance, other factors (such as the complexity of the search loop and the overhead is computing the hash address) also affect the running time when Algorithm C is programmed for a real computer. For completeness, we optimize the running time of assembly language versions of coalesced hashing.

We choose to program in assembly language rather than in some high-level language like Fortran, PL/I, or Pascal, in order to achieve *maximum possible efficiency*. Top efficiency is important in large-scale applications of hashing, but it can also be achieved in smaller systems with little extra effort, because hashing algorithms are so short that implementing them (even in assembly language) is easy. We use a hypothetical language based on Knuth's MIX [6] because its features are similar to most well-known machines and its inherent simplicity allows us to write programs in clear and concise form.

Program C below is a MIX-like implementation of Algorithm C. Liberties have been taken with the language for purposes of clarity; the actual MIX code appears in [10]. The program is written in a five-column format: the first column gives the line numbers, the second column lists the instruction labels, the third column contains the assembly language instructions, the fourth column counts the number of times the instructions are executed, and the last column is for comments that explain what the instructions do. The syntax of the commands should be clear to those familiar with assembly language programming. The four memory registers

used in Program C are named rA , rX , rI , and rJ . The reference $KEY(I)$ denotes the contents of the memory location whose address is the value of KEY plus the contents of rI . (This is $KEY[i]$ in the notation of Algorithm C.)

Program C (*Coalesced hashing search and insertion*). This program follows the conventions of Algorithm C, except that the *EMPTY* field is implicit in the *LINK*

field: empty slots are marked by a -1 in the *LINK* field of that slot. Null links are denoted by a 0 in the *LINK* field. The variable R and the key K are stored in memory locations R and K . Registers rI and rA are used to store the values of I and K . Register rJ stores either the value of $LINK[I]$ or R . The instruction labels *SUCCESS* and *OVERFLOW* are for exiting and are assumed to lie somewhere outside this code.

| | | | | | |
|----|-------|-----|-------------|------------|---|
| 01 | START | LD | X, K | 1 | Step C1. Load rX with K . |
| 02 | | ENT | A, 0 | 1 | Enter 0 into rA . |
| 03 | | DIV | =M= | 1 | $rA \leftarrow \lfloor K/M \rfloor$, $rX \leftarrow K \bmod M$. |
| 04 | | ENT | I, X | 1 | Enter rX into rI . |
| 05 | | INC | I, 1 | 1 | Increment rI by 1. |
| 06 | | LD | A, K | 1 | Load rA with K . |
| 07 | | LD | J, LINK(I) | 1 | Step C2. Load rJ with $LINK[I]$. |
| 08 | | JN | J, STEP6 | 1 | Jump to STEP6 if $LINK[I] < 0$. |
| 09 | | CMP | A, KEY(I) | A | Step C3. Compare K with $KEY[I]$. |
| 10 | | JE | SUCCESS | A | Exit (successfully) if $K = KEY[I]$. |
| 11 | | JZ | J, STEP5 | A - S1 | Jump to STEP5 if $LINK[I] = 0$. |
| 12 | STEP4 | ENT | I, J | C - 1 | Step C4. Enter rJ into rI . |
| 13 | | CMP | A, KEY(I) | C - 1 | Step C3. Compare K with $KEY[I]$. |
| 14 | | JE | SUCCESS | C - 1 | Exit (successfully) if $K = KEY[I]$. |
| 15 | | LD | J, LINK(I) | C - 1 - S2 | Load rJ with $LINK[I]$. |
| 16 | | JNZ | J, STEP4 | C - 1 - S2 | Jump to STEP4 if $LINK[I] \neq 0$. |
| 17 | STEP5 | LD | J, R | A - S | Step C5. Load rJ with R . |
| 18 | | DEC | J, 1 | T | Decrement R by 1. |
| 19 | | LD | X, LINK(J) | T | Load rX with $LINK[R]$. |
| 20 | | JNN | X, -2 | T | Go back two steps if $LINK[R] \geq 0$. |
| 21 | | JZ | J, OVERFLOW | A - S | Exit (with overflow) if $R = 0$. |
| 22 | | ST | J, LINK(I) | A - S | Store R in $LINK[I]$. |
| 23 | | ENT | I, J | A - S | Enter rJ into rI . |
| 24 | | ST | J, R | A - S | Update R in memory. |
| 25 | STEP6 | ST | 0, LINK(I) | 1 - S | Step C6. Store 0 in $LINK[I]$. |
| 26 | | ST | A, KEY(I) | 1 - S | Store K in $KEY[I]$. ■ |

The execution time is measured in *MIX units of time*, which we denote u . The number of time units required by an instruction is equal to the number of memory references (including the reference to the instruction itself). Hence, the LD, ST, and CMP instructions each take two units of time, while ENT, INC, DEC, and the jump instructions require only one time unit. The division operation used to compute the hash address is an exception to this rule; it takes $14u$ to execute.

The running time of a MIX program is the weighted sum

$$\sum_{\text{each instruction in the program}} \left(\begin{array}{c} \# \text{ times} \\ \text{the instruction} \\ \text{is executed} \end{array} \right) \left(\begin{array}{c} \# \text{ time units} \\ \text{required by} \\ \text{the instruction} \end{array} \right) \quad (1)$$

This is a somewhat simplistic model, since it does not make use of cache or buffered memory for fast access of frequently used data, and since it ignores any intervention by the operating system. But it places all hashing algorithms on an equal footing and gives a good indication of relative merit.

The fourth column of Program C expresses the number of times each instruction is executed in terms of the quantities

C = number of probes per search.

$A = 1$ if the initial probe found an occupied slot, 0 otherwise.

$S = 1$ if successful, 0 if unsuccessful.

T = number of slots probed while looking for an empty space.

We further decompose S into $S1 + S2$, where $S1 = 1$ if the search is successful on the first probe, and $S1 = 0$ otherwise. By formula (1), the total running time of the searching phase is

$$(7C + 4A + 17 - 3S + 2S1)u \quad (2)$$

and the insertion of a new record after an unsuccessful search (when $S = 0$) takes an additional $(8A + 4T + 4)u$. The average running time is the expected value of (2), assuming that all insertions and searches are random. The formula can be obtained by replacing the variables in Eq. (2) with their expected values.

4. Tuning β to Obtain Optimum Performance

The purpose of the analysis in [10, 11, 13] is to show how the average-case performance of the coalesced hashing method varies as a function of the address factor $\beta = M/M'$ and the load factor $\alpha = N/M'$. In this section, for each fixed value of α , we make use of those results in order to "tune" our choice of β and speed up the search times. Our two measures of performance are the expected number of probes per search and the average running time of assembly language versions. In the latter case, we study a MIX implementation in detail, and then show how to apply what we learn to other assembly languages.

Unfortunately, there is no single choice of β that yields best results: the optimum choice β_{OPT} is a function of the load factor α and it is even different for unsuccessful and successful searches. The section concludes with practical tips on how to initialize β . In particular, we shall see that the choice $\beta \approx 0.86$ works well in most situations.

4.1 Number of Probes Per Search

For each fixed value of α , we want to find the values β_{OPT} that minimize the expected number of search probes in unsuccessful and successful searches. Formulas (A1) and (A2) in the Appendix express the average number of probes per search as a function of three variables: the load factor $\alpha = N/M'$, the address factor $\beta = M/M'$, and a new variable $\lambda = L/M$, where L is the expected number of inserted records needed to make the cellar become full. The variables β and λ are related by the formula

$$e^{-\lambda} + \lambda = \frac{1}{\beta} \quad (3)$$

Formulas (A1) and (A2) each have two cases, " $\alpha \leq \lambda\beta$ " and " $\alpha \geq \lambda\beta$," which have the following intuitive meanings: The condition $\alpha < \lambda\beta$ means that with high probability not enough records have been inserted to fill up the cellar, while the condition $\alpha > \lambda\beta$ means that enough records have been inserted to make the cellar almost surely full.

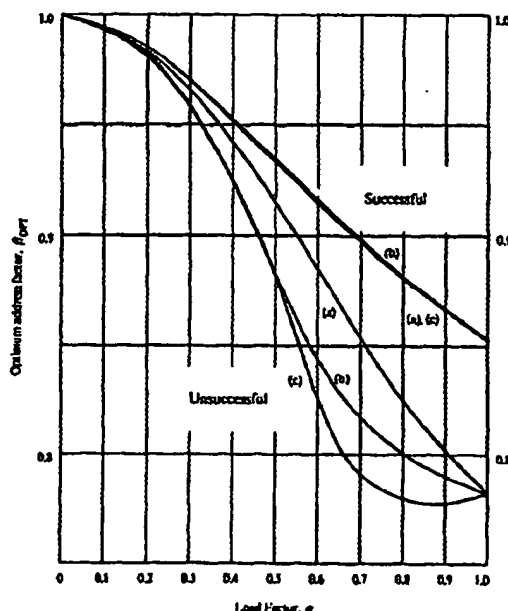
The optimum address factor β_{OPT} is always located somewhere in the " $\alpha \geq \lambda\beta$ " region, as shown in the Appendix. The rest of the optimization procedure is a straightforward application of differential calculus. First, we substitute Eq. (3) into the " $\alpha \geq \lambda\beta$ " cases of the formulas for the expected number of probes per search in order to express them in terms of only the two variables α and λ . For each nonzero fixed value of α , the formulas are convex w.r.t. λ and have unique minima. We minimize them by setting their derivatives equal to 0. Numerical analysis techniques are used to solve the resulting equations and to get the optimum values of λ for several different values of α . Then we reapply Eq. (3) to express the optimum points in terms of β . The results are graphed in Fig. 2(a), using spline interpolation to fill in the gaps.

4.2 MIX Running Times

Optimizing the MIX execution times could be tricky, in general, because the formulas might have local as well as global minima. Then when we set the derivatives equal to 0 in order to find β_{OPT} , there might be several roots to the resulting equations. The crucial fact that lets us apply the same optimization techniques we used above for the number of probes is that the formulas for the MIX running times are *well-behaved*, as shown in the Appendix. By that we mean that each formula is minimized at a *unique* β_{OPT} , which occurs either at the endpoint $\alpha = \lambda\beta$ or at the unique point in the " $\alpha \geq \lambda\beta$ " region where the derivative w.r.t. β is 0.

The optimization procedure is the same as before. The expected values of formulas (A4) and (A5), which give the MIX running times for unsuccessful and successful searches, are functions of the three variables α , β , and λ . We substitute Eq. (3) into the expected running times in order to express β in terms of λ . For several different load factors α and for each type of search, we find the value of λ that minimizes the formula, and then we retranslate this value via Eq. (3) to get β_{OPT} . Figure 2(b) graphs these optimum values β_{OPT} as a function of α ; spline interpolation was used to fill in the gaps. As in the previous section, the formulas for the average unsuccessful and successful search times yield different optimum address factors. For the successful search case, notice how closely β_{OPT} agrees with the corresponding values that minimize the expected number of probes.

Fig. 2. The values β_{OPT} that optimize search performance for the following three measures: (a) the expected number of probes per search, (b) the expected running time of Program C, and (c) the expected assembly language running time for large keys.



4.3 Applying the Results to Other Implementations

Our MIX analysis suggests two important principles to be used in finding β_{OPT} for a particular implementation of coalesced hashing. First, the formulas for the expected number of times each instruction in the program is executed (which are expressed for Program C in terms of $C, A, S, S1, S2$, and T) may have the two cases, " $\alpha \leq \lambda\beta$ " and " $\alpha \geq \lambda\beta$," but probably not more.

Second, the same optimization process as above can be used to find β_{OPT} , because the formulas for the running times should be well-behaved for the following reason: The main difference between Program C and another implementation is likely to be the relative time it takes to process each key. (The keys are assumed to be very small in the MIX version.) Thus, the unsuccessful search time for another implementation might be approximately

$$[(2\kappa + 5)C + (2\kappa + 2)A + (-2\kappa + 19)u'] \quad (4)$$

where u' is the standard unit of time on the other computer and κ is how many times longer it takes to process a key (multiplied by u/u'). Successful search times would be about

$$(2\kappa + 5)C + 18 + 2S1)u' \quad (5)$$

Formulas (4) and (5) were calculated by increasing the execution times of the key-processing steps 9 and 13 in Program C by a factor of κ . (See formulas (A4) and (A5) for the $\kappa = 1$ case.) We ignore the extra time it takes to load the larger key and to compute the hash function, since that does not affect the optimization.

The role of C in formula (4) is less prevalent than in (A4) as κ gets large: the ratio of the coefficients of C and A decreases from $7/4$ in (A4) and approaches the limit $2/2 = 1$ in formula (4). Even in this extreme case, however, computer calculations show that the formula for the average running time is well-behaved. The values of β_{OPT} that minimize formula (4) when κ is large are graphed in Fig. 2(c).

For successful searches, however, the value of C more strongly dominates the running times for larger values of κ , so the limiting values of β_{OPT} in Fig. 2(c) coincide with the ones that minimize the expected number of probes per search in Fig. 2(a). Figure 2(b) shows that the approximation is close even for the case $\kappa = 1$, which is Program C.

4.4 How to Choose β

It is important to remember that the address region size $M = \lceil \beta M' \rceil$ must be initialized when the hash table is empty and cannot change thereafter. Unfortunately, the last two sections show that each different load factor α requires a different optimum address factor β_{OPT} ; in fact, the values of β_{OPT} differ for unsuccessful and successful searches. This means that optimizing the average unsuccessful (or successful) search time for a certain load factor α will lead to suboptimum performance when the load factor is not equal to α .

One strategy is to pick $\beta \approx 0.782$, which minimizes the expected number of probes per unsuccessful search as well as the average MIX unsuccessful search time when the table is full (i.e., load factor $\alpha = 1$), as indicated in Fig. 2. This choice of β yields the best absolute bound on search performance, because when the table is full, search times are greatest and unsuccessful searches average slightly longer than successful ones. Regardless of the load factor, the expected number of probes per search would be at most 1.79, and the average MIX searching time would be bounded by $33.52u$.

Another strategy is to pick some compromise address factor that leads to good overall performance for a large range of load factors. A reasonable choice is $\beta = 0.86$; then the unsuccessful searches are optimized (over all other values of β) when the load factor is ≈ 0.68 (number of probes) and ≈ 0.56 (MIX), and the successful search performance is optimized at load factors ≈ 0.94 (number of probes) and ≈ 0.95 (MIX).

Figures 3 through 6 graph the expected search performance of coalesced hashing as a function of α for both types of searches (unsuccessful and successful) and for both measures of performance (number of probes and MIX running time). The C_1 curve corresponds to standard coalesced hashing (i.e., $\beta = 1$); the $C_{0.86}$ line is our compromise choice $\beta = 0.86$; and the dashed line C_{OPT} represents the best possible search performance that could be achieved by tuning (in which β is optimized for each load factor).

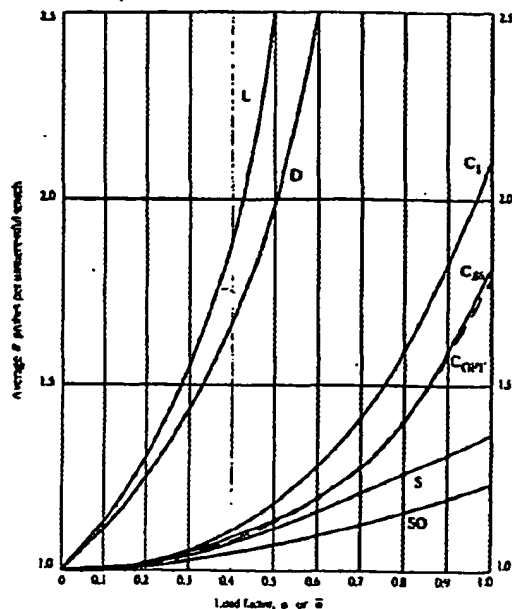
Notice that the value $\beta = 0.86$ yields near-optimum search times once the table gets half-full, so this compromise offers a viable strategy. Of course, if some prior knowledge about the types and frequencies of the searches were available, we could tailor our choice of β to meet those specific needs.

5. Comparisons

In this section, we compare the searching times of the coalesced hashing method with those from a representative collection of hashing schemes: standard coalesced hashing (C_1), separate chaining (S), separate chaining with ordered chains (SO), linear probing (L), and double hashing (D). Implementations of the methods are given in [10].

These methods were chosen because they are the most well-known and since they each have implementations similar to that of Algorithm C. Our comparisons are based both on the expected number of probes per search as well as on the average MIX running time. Coalesced hashing performs better than the other methods. The differences are not so dramatic with the MIX search times as with the number of probes per search, due to the large overhead in computing the hash address. However, if the keys were larger and comparisons took longer, the relative MIX savings would closely approximate the savings in number of probes.

Fig. 3. The average number of probes per unsuccessful search, as M and $M' \rightarrow \infty$, for coalesced hashing (C_1 , C_{cas} , C_{opt} for $\beta = 1, 0.86, \beta_{\text{opt}}$), separate chaining (S), separate chaining with ordered chains (SO), linear probing (L), and double hashing (D).



5.1 Standard Coalesced Hashing (C_1)

Standard coalesced hashing is the special case of coalesced hashing for which $\beta = 1$ and there is no cellar. This is obviously the most realistic comparison that can be made, because except for the initialization of the address region size, standard coalesced hashing and

Fig. 5. The average MIX execution time per unsuccessful search, as $M' \rightarrow \infty$, for coalesced hashing (C_1 , C_{cas} , C_{opt} for $\beta = 1, 0.86, \beta_{\text{opt}}$), separate chaining (S), separate chaining with ordered chains (SO), linear probing (L), and double hashing (D).

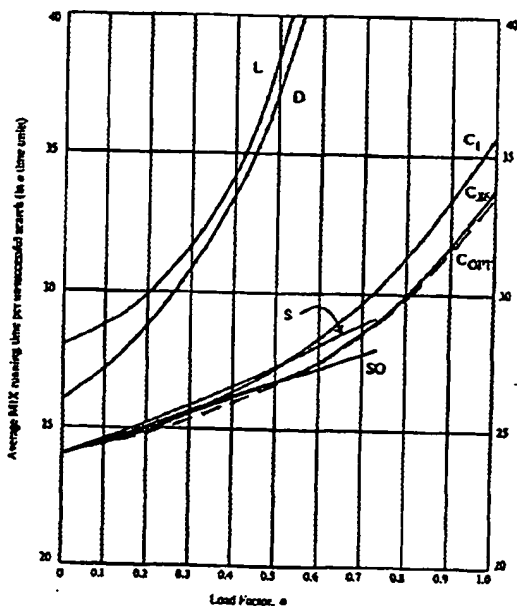
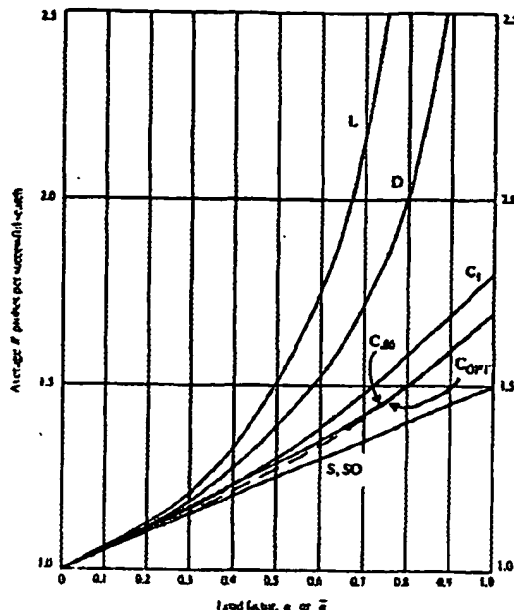
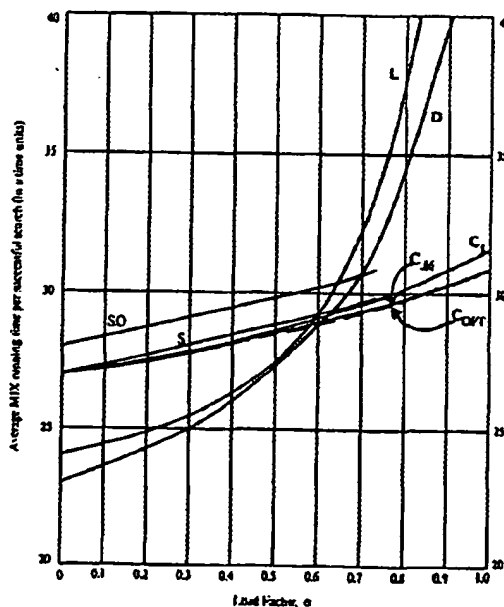


Fig. 4. The average number of probes per successful search, as M and $M' \rightarrow \infty$, for coalesced hashing (C_1 , C_{cas} , C_{opt} for $\beta = 1, 0.86, \beta_{\text{opt}}$), separate chaining (S), separate chaining with ordered chains (SO), linear probing (L), and double hashing (D).



"tuned" coalesced hashing are identical. Figures 3 and 4 show that the savings in number of probes per search can be as much as 14 percent (unsuccessful) and 6 percent (successful). In Figs. 5 and 6, the corresponding savings in MIX searching time is 6 percent (unsuccessful) and 2 percent (successful).

Fig. 6. The average MIX execution time per successful search, as $M' \rightarrow \infty$, for coalesced hashing (C_1 , C_{cas} , C_{opt} for $\beta = 1, 0.86, \beta_{\text{opt}}$), separate chaining (S), separate chaining with ordered chains (SO), linear probing (L), and double hashing (D).



5.2 Separate (or Direct) Chaining (S)

The separate chaining method is given an unfair advantage in Figs. 3 and 4: the number of probes per search is graphed as a function of $\bar{\alpha} = N/M$ rather than $\alpha = N/M'$ and does not take into account the number of auxiliary slots used to store colliders. In order to make the comparison fair, we must adjust the load factor accordingly.

Separate chaining implementations are designed often to accommodate about $N = M$ records; an average of $M(1 - 1/M') \approx M/e$ auxiliary slots are needed to store the colliders. The total table size is thus $M' \approx M + M/e$. Solving backwards for M , we get $M \approx 0.731M'$. In other words, we may consider separate chaining to be the special case of coalesced hashing for which $\beta \approx 0.731$, except that no more records can be inserted once the cellar overflows. Hence, the adjusted load factor is $\alpha \approx 0.731\bar{\alpha}$, and overflow occurs when there are around $N = M \approx 0.731M'$ inserted records. (This is a reasonable space/time compromise: if we make M smaller, then more records can usually be stored before overflow occurs, but the average search times blow up; if we increase M to get better search times, then overflow occurs much sooner, and many slots are wasted.)

If we adjust the load factors in Figs. 3 and 4 in this way, Algorithm C generates better search statistics: the expected number of probes per search for separate chaining is ≈ 1.37 (unsuccessful) and ≈ 1.5 (successful) when the load factor $\bar{\alpha}$ is 1, while that for coalesced hashing is ≈ 1.32 (unsuccessful) and ≈ 1.44 (successful) when the load factor $\alpha = \beta\bar{\alpha}$ is equal to 0.731.

The graphs in Figs. 5 and 6 already reflect this load factor adjustment. In fact, the MIX implementation of separate chaining (Program S in [10]) is identical to Program C, except that β is initialized to 0.731 and overflow is signaled automatically when the cellar runs out of empty slots. Program C is slightly quicker in MIX execution time than Program S, but more importantly, the coalesced hashing implementation is more space efficient: Program S usually overflows when $\alpha \approx 0.731$, while Program C can always obtain full storage utilization $\alpha = 1$. This confirms our intuition that coalesced hashing can accommodate more records than the separate chaining method and still outperform separate chaining before that method overflows.

5.3 Separate Chaining with Ordered Chains (SO)

This method is a variation of separate chaining in which the chains are kept ordered by key value. The expected number of probes per successful search does not change, but unsuccessful searches are slightly quicker, because only about half the chain needs to be searched, on the average.

Our remarks about adjusting the load factor in Figs. 3 and 4 also apply to method SO. But even after that is done, the average number of probes per unsuccessful search as well as the expected MIX unsuccessful search time is slightly better for this method than for coalesced hashing. However, as Fig. 6 illustrates, the average suc-

cessful search time of Program SO is worse than Program C's, and in real-life situations, the difference is likely to be more apparent, because records that are inserted first tend to be looked up more often and should be kept near the beginning of the chain, not rearranged.

Method SO has the same storage limitations as the separate chaining scheme (i.e., the table usually overflows when $N \approx M \approx 0.731M'$), whereas coalesced hashing can obtain full storage utilization.

5.4 Linear Probing (L) and Double Hashing (D)

When searching for a record with key K , the linear probing method first checks location $hash(K)$, and if another record is already there, it steps cyclically through the table, starting at location $hash(K)$, until the record is found (successful search) or an empty slot is reached (unsuccessful search). Insertions are done by placing the record into the empty slot that terminated the unsuccessful search. Double hashing generalizes this by letting the cyclic step size be a function of K .

We have to adjust the load factor in the *opposite* direction when we compare Algorithm C with methods L and D, because the latter do not require *LINK* fields. For example, if we suppose that the *LINK* field comprises $\frac{1}{4}$ of the total record size in a coalesced hashing implementation, then the search statistics in Figs. 3 and 4 for Algorithm C with load factor α should be compared against those for linear probing and double hashing with load factor $(\frac{3}{4})\alpha$. In this case, the average number of probes per search is still better for coalesced hashing.

However, the *LINK* field is often much smaller than the rest of the record, and sometimes it can be included in the table at virtually no extra cost. The MIX implementation Program C in [10] assumes that the MIX field can be squeezed into the record without need of extra storage space. Figures 5 and 6, therefore, require no load factor adjustment.

To balance matters, the MIX implementations of linear probing and double hashing, which are given in [10] and [7], contain two code optimizations. First, since *LINK* fields are not used in methods L and D, we no longer need 0 to denote a null *LINK*, and we can renumber the table slots from 0 to $M' - 1$; the hash function now returns a value between 0 and $M' - 1$. This makes the hash address computation faster by 1u, because the instruction INC I, 1 can be eliminated. Second, the empty slots are denoted by the value 0 in order to make the comparisons in the inner loop as fast as possible. This means that records are not allowed to have a key value of 0. The final results are graphed in Figs. 5 and 6. Coalesced hashing clearly dominates when the load factor is greater than 0.6.

6. Deletions

It is often useful in hashing applications to be able to delete records when they no longer logically belong to the set of objects being represented in the hash table. For

example, in an airlines reservations system, passenger records are often expunged soon after the flight has taken place.

One possible deletion strategy often used for linear probing and double hashing is to include a special one-bit *DELETED* field in each record that says whether or not the record has been deleted. The search algorithm must be modified to treat each "deleted" table slot as if it were occupied by a null record, even though the entire record is still there. This is especially desirable when there are pointers to the records from *outside* the table.

If there are no such external pointers to worry about, the "deleted" table slots can be reused for later insertions: Whenever an empty slot is needed in step C5 of Algorithm C, the record is inserted into the first "deleted" slot encountered during the unsuccessful search; if there is no such slot, an empty slot is allocated in the usual way. However, a certain percentage of the "deleted" slots probably will remain unused, thus preventing full storage utilization. Also, insertions and deletions over a prolonged period would cause the expected search times to approximate those for a full table, regardless of the number of undeleted records, because the "deleted" records make the searches longer.

If we are willing to spend a little extra time per deletion, we can do without the *DELETED* field by relocating some of the records that follow in the chain. The basic idea is this: First, we find the record we want to delete, mark its table slot empty, and set the *LINK* field of its predecessor (if any) to the null value 0. Then we use Algorithm C to reinsert each record in the remainder of the chain, but whenever an empty slot is needed in step C5, we use the position that the record already occupies.

This method can be illustrated by deleting AL from location 10 in Fig. 7(a); the end result is pictured in Fig. 7(b). The first step is to create a hole in position 10 where AL was, and to set AUDREY's *LINK* field to 0. Then we process the remainder of the chain. The next record

TOOTIE rehashes to the hole in location 10, so TOOTIE moves up to plug the hole, leaving a new hole in position 9. Next, DONNA collides with AUDREY during rehashing, so DONNA remains in slot 8 and is linked to AUDREY. Then MARK also collides with AUDREY; we leave MARK in position 7 and link it to DONNA, which was formerly at the end of AUDREY's hash chain. The record JEFF rehashes to the hole in slot 9, so we move it up to plug the hole, and a new hole appears in position 6. Finally, DAVE rehashes to position 9 and joins JEFF's chain.

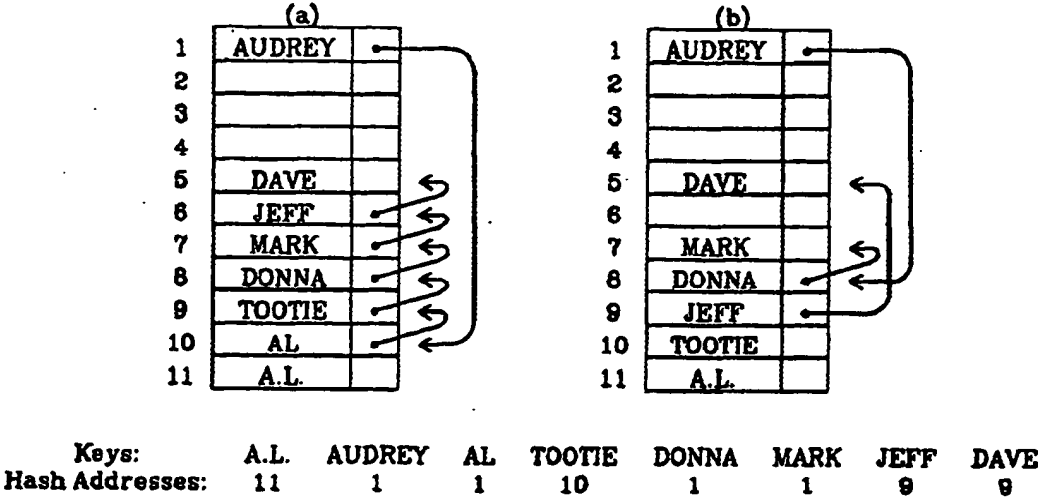
Location 6 is the current hole position when the deletion algorithm terminates, so we set *EMPTY*[6] ← true and return it to the pool of empty slots. However, the value of *R* in Algorithm C is already 5, so step C5 will never try to reuse location 6 when an empty slot is needed.

We can solve this problem by using an *available-space list* in step C5 rather than the variable *R*; the list must be doubly linked so that a slot can be removed quickly from the list in step C6. The available-space list does not require any extra space per table slot, since we can use the *KEY* and *LINK* fields of the empty slots for the two pointer fields. (The *KEY* field is much larger than the *LINK* field in typical implementations.) For clarity, we rename the two pointer fields *NEXT* and *PREV*. Slot 0 in the table acts as the dummy start of the available-space list, so *NEXT*[0] points to the first actual slot in the list and *PREV*[0] points to the last. Before any records are inserted into the table, the following extra initializations must be made: *NEXT*[0] ← *M'*, *PREV*[*M'*] ← 0; and *NEXT*[*i*] ← *i* - 1 and *PREV*[*i* - 1] ← *i*, for 1 ≤ *i* ≤ *M'*. We replace steps C5 and C6 by

C5. [Find empty slot.] (The search for *K* in the chain was unsuccessful, so we will try to find an empty table slot to store *K*.) If the table is already full (i.e., *NEXT*[0] = 0), the algorithm terminates with *overflow*. Otherwise, set *LINK*[*i*] ← *NEXT*[0] and *i* ← *NEXT*[0].

C6. [Insert new record.] Remove the *i*th slot from the

Fig. 7. (a) Inserting the eight records; (b) Inserting all the records *except* AL.



**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☒ **BLACK BORDERS**
- ☐ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**
- ☐ **FADED TEXT OR DRAWING**
- ☒ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**
- ☐ **SKEWED/SLANTED IMAGES**
- ☐ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**
- ☐ **GRAY SCALE DOCUMENTS**
- ☐ **LINES OR MARKS ON ORIGINAL DOCUMENT**
- ☐ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**
- ☐ **OTHER:** _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.